

precisely

Spectrum Data Quality

Data Quality User Guide

Version 2020.1.0



Table of Contents

1 - Getting Started

Introduction to Data Quality.....	5
-----------------------------------	---

2 - Parsing

Introduction to Parsing.....	8
Defining Domain-Independent Parsing Grammars in Dataflows.....	9
Culture-Specific Parsing.....	10
Analyzing Parsing Results.....	36
Parsing Personal Names.....	40
Dataflow Templates for Parsing.....	41

3 - Standardization

Standardizing Terms.....	56
Standardizing Personal Names.....	57
Templates for Standardization.....	59

4 - Matching

Matching Terminology.....	62
Standard Fields	63
Techniques for Defining Match Keys.....	65
Match Rules.....	68
Matching Records from a Single Source.....	80
Matching Records from One Source to Another Source.....	86
Matching Records Between and Within Sources.....	92
Matching Records Against a Database.....	98
Matching Records Using Multiple Match Rules.....	100
Creating a Universal Matching Service.....	104
Using an Express Match Key.....	107

Analyzing Match Results.....	111
Dataflow Templates for Matching.....	126

5 - Deduplication

Filtering Out Duplicate Records.....	135
Creating a Best of Breed Record.....	139

6 - Exception Records

Designing a Dataflow to Handle Exceptions.....	145
Designing a Dataflow for Real-Time Revalidation.....	146

7 - Lookup Tables

Introduction to Lookup Tables.....	150
Data Normalization Tables.....	150
Universal Name Tables.....	155
Viewing the Contents of a Lookup Table.....	157
Adding a Term to a Lookup Table.....	158
Removing a Term from a Lookup Table.....	158
Modifying the Standardized Form of a Term.....	159
Reverting Table Customizations.....	159
Creating a Lookup Table.....	160
Importing Data.....	160

8 - Data Stewardship Portal

Introduction.....	164
Accessing the Data Stewardship Portal.....	164
The Dashboard page.....	166
The Data Stewardship Portal Editor.....	167
The Manage page.....	193
The Data Quality Page.....	198

How to Video - Approval Flow in Data	
Stewardship.....	199
Introduction.....	199
Accessing the Data Stewardship Portal.....	200
The Dashboard page.....	201
The Data Stewardship Portal Editor.....	202
The Manage page.....	229
The Data Quality Page.....	234
How to Video - Approval Flow in Data	
Stewardship.....	235
Data Stewardship Settings.....	235

9 - Accessing Stages in Enterprise Designer

Advanced Matching.....	246
Data Stewardship.....	310
Data Normalization.....	325
Universal Name.....	344

10 - Accessing Stages through Flow Designer

Advanced Matching.....	380
------------------------	-----

11 - Performance Tuning

Performance Tuning Recommendations.....	427
---	-----

12 - ISO Country Codes and Module Support

ISO Country Codes and Module Support.....	431
---	-----

1 - Getting Started

In this section

Introduction to Data Quality.....5



Introduction to Data Quality

Data Quality involves ensuring the accuracy, timeliness, completeness, and consistency of the data used by an organization so that the data is fit for use. Spectrum Technology Platform supports data quality initiatives by providing the following capabilities.

Parsing

Parsing is the process of analyzing a sequence of input characters in a field and breaking it up into multiple fields. For example, you might have a field called Name which contains the value "John A. Smith" and through parsing, you can break it up so that you have a FirstName field containing "John", a MiddleName field containing "A" and a LastName field containing "Smith."

Standardization

Standardization takes data of the same type and puts it in the same format. Some types of data that may be standardized include telephone numbers, dates, names, addresses, and identification numbers. For example, telephone numbers can be formatted to eliminate non-numeric characters such as parentheses, periods, or dashes.

You should standardize your data before performing matching or deduplication activities since standardized data will be more accurately matched than data that is inconsistently formatted.

Matching

Matching is the process of identifying records that are related to each other in some way that is significant for your purposes. For example, if you are trying to eliminate redundant information from your customer data, you may want to identify duplicate records for the same customer; or, if you are trying to eliminate duplicate marketing pieces going to the same address, you may want to identify records of customers that live in the same household.

Deduplication

Deduplication identifies records that represent one entity but for one reason or another were entered into the system multiple times, sometimes with slightly different data. For example, your system may contain vendor information from different departments in your organization, with each department using a different vendor ID for the same vendor. Using Spectrum Technology Platform you can consolidate these records into a single record for each vendor.

Review of Exception Records

In some cases you may have data that cannot be confidently processed automatically and that must be reviewed by a knowledgeable data steward. Some examples of records that may require manual review include:

- Address verification failures
- Geocoding failures
- Low-confidence matches
- Merge/consolidation decisions

The Data Stewardship Portal contains features that allow you to identify and resolve exception records.

2 - Parsing

In this section

Introduction to Parsing.....	8
Defining Domain-Independent Parsing Grammars in Dataflows.....	9
Culture-Specific Parsing.....	10
Analyzing Parsing Results.....	36
Parsing Personal Names.....	40
Dataflow Templates for Parsing.....	41



Introduction to Parsing

Parsing is the process of analyzing a sequence of input characters in a field and breaking it up into multiple fields. For example, you might have a field called Name which contains the value "John A. Smith" and through parsing, you can break it up so that you have a FirstName field containing "John", a MiddleName field containing "A" and a LastName field containing "Smith."

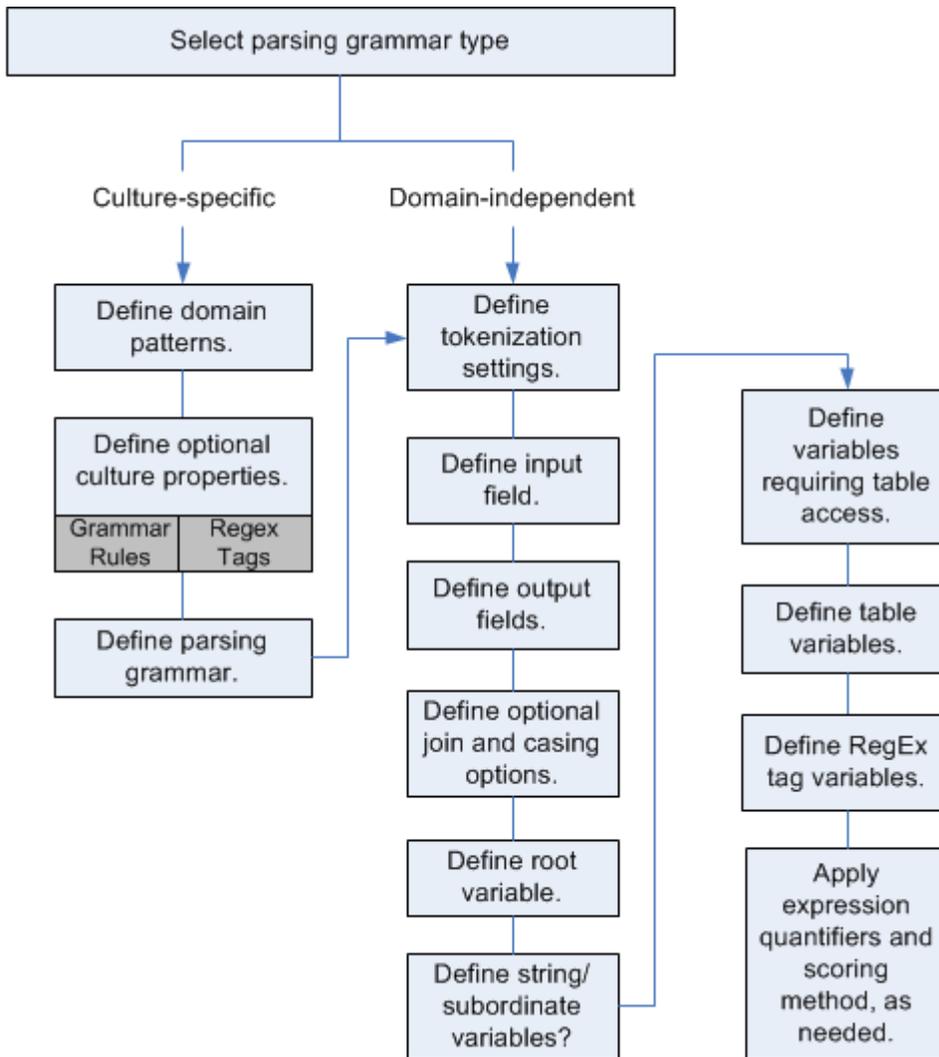
To create a dataflow that parses, use the Open Parser stage. Open Parser allows you to write parsing rules called *grammars*. A grammar is a set of expressions that map a sequence of characters to a set of named entities called *domain patterns*. A domain pattern is a sequence of one or more tokens in your input data that you want to represent as a data structure, such as name, address, or account numbers. A domain pattern can consist of any number of tokens that can be parsed from your input data. A domain pattern is represented in the parsing grammar as the <root> expression. Input data often contains such tokens in hard-to-use or mixed formats. For example:

- Your input data contains names in a single field that you want to separate into given name and family name.
- Your input data contains addresses from several cultures and you want to extract address data for a specific culture only.
- Your input data includes free-form text that contains embedded email addresses and you want to extract email addresses and match them up with personal data and store them in a database.

There are two kinds of grammars: *culture specific* and *domain independent*. A culture-specific parsing grammar is associated with a culture and/or language (such as English, Canadian English, Spanish, Mexican Spanish, and so on) and a particular type of data (phone numbers, personal names, and so on). When an Open Parser stage is configured to perform culture-specific parsing, each culture's parsing grammar is applied to each record. The grammar with the best parser score (or the first one to have a score of 100) is the one whose results are returned. Alternatively, culture-specific parsing grammars can use the value in the input record's CultureCode field and process the data according to the culture settings contained in the culture's parsing grammar. Culture-specific parsing grammars can inherit properties from a parent. A domain-independent parsing grammar is not associated with either a language or a particular type of data. Domain-independent parsing grammars do not inherit properties from a parent and ignore any CultureCode information in the input data.

Open Parser analyzes a sequence of characters in input fields and categorizes them into a sequence of tokens through a process called *tokenization*. Tokenization is the process of delimiting and classifying sections of a string of input characters into a set of tokens based on separator characters (also called tokenizing characters), such as space, hyphen, and others. The tokens are then placed into output fields you specify.

The following diagram illustrates the process of creating a parsing grammar:



Defining Domain-Independent Parsing Grammars in Dataflows

To define domain-independent parsing grammars in a dataflow:

1. In Enterprise Designer, add an Open Parser stage to your dataflow.
2. Double-click the Open parser stage on the canvas.
3. Click **Define Domain Independent Grammar** on the **Rules** tab.

4. Use the Grammar Editor to create the grammar rules. You can type commands and variables into the text box or use the commands provided in the **Commands** tab. For more information, see [Grammars](#) on page 27.
5. To cut, copy, paste, and find and replace text strings in your parsing grammar, right-click in the Grammar Editor and select the appropriate command.
6. To check the parsing grammar you have created, click **Validate**.

The validate feature lists any errors in your grammar syntax, including the line and column where the error occurs, a description of the error, and the command name or value that is causing the error.

7. Click the **Preview** tab to test the parsing grammar.
8. When you are finished creating your parsing grammar, click **OK**.

Culture-Specific Parsing

Defining a Culture-Specific Parsing Grammar

A culture-specific parsing grammar allows you to specify different parsing rules for different languages and cultures. This allows you to parse data from different countries in a single Open Parser stage, for example phone numbers from the United States and phone numbers from the United Kingdom. By default, each input record is parsed using each culture's parsing grammar, in the order specified in the Open Parser stage. You can also add a CultureCode field to the input records if you want a specific culture's parsing grammar to be used for that record. For more information, see [Assigning a Parsing Culture to a Record](#) on page 12.

Note: If you want to create a domain-independent parsing grammar, see [Defining Domain-Independent Parsing Grammars in Dataflows](#) on page 9.

1. In Enterprise Designer, go to **Tools > Open Parser Domain Editor**.
2. Click the **Domains** tab.
3. Click **Add**.
4. Type a domain name in the **Name** field.
5. Type a description of the domain name in the **Description** field.
6. If you want to create a new, empty domain, click **OK**. If you want to create a new domain based on another domain, do the following:
 - a) Select **Use another domain as a template** if you want to create a new domain based on another domain.

- b) Select a domain from the list. When you click **OK** in the next step, the new domain will be created. The new domain will contain all of the culture-specific parsing grammars defined in the domain template that you selected.
 - c) Click **OK**.
7. Define the parsing grammar for the global culture. The global culture is the default culture and is used to parse records that have a culture for which no culture-specific parsing grammar has been defined.
 - a) On the **Grammars** tab, select the new domain you created.
 - b) If you created a domain from a template, there may be cultures already listed.
 - If there are cultures listed, select **Global Culture** then click **Edit**.
 - If there are no cultures listed, click **Add**, select **Global Culture** then click **OK**.
 - c) On the **Grammar** tab, write the parsing grammar for the global culture. You can use the **Commands**, **Grammar Rules**, and **RegEx Tags** tabs to insert predefined parsing grammar elements. To enter a predefined element, place the cursor where you want to insert the element then double-click the element you want to add.

The **Commands** tab displays parsing commands. For information about the commands available, see [Grammars](#) on page 27.

The **Grammar Rules** tab displays grammar rules that you create in the Culture Properties dialog box. For more information about creating grammar rules, see [Defining a Culture's Grammar Rules](#) on page 31.

The **RegEx Tags** tab displays RegEx tags that you create in the Culture Properties dialog box. For more information about creating RegEx tags, see [Defining Culture RegEx Tags](#) on page 32.
 - d) To check the grammar syntax you have created, click **Validate**. The parsing grammar validation feature displays any errors in your grammar syntax and includes the error encountered, the line and column where the error occurs, and the command, grammar rule, or RegEx tag where the error occurs.
 - e) To test the results of your grammar with sample data, click the **Preview** tab. Under **Input Data**, enter sample data you want to parse. Enter one record per row. Then, click the **Preview** button. The parsed output fields display in the **Results** grid. For information about the output fields, see [Output](#) on page 333. For information about trace, see [Tracing Final Parsing Results](#) on page 36. If your results are not what you expected, click the **Grammars** tab and continue editing the parsing grammar and testing representative input data until the parsing grammar produces the expected results.
 - f) Click **OK** when you are done defining the parsing grammar for the global culture.
8. Define a culture-specific grammar for each culture you want. To add culture-specific grammars, click **Add** and define the grammar using the same steps as for the global culture. Repeat as needed to add as many cultures as you need.
9. When you are done adding culture-specific parsing grammars, click **OK**.

The domain and cultures you have created can now be used in the Open Parser stage to perform parsing.

Assigning a Parsing Culture to a Record

When you configure an Open Parser stage to use culture-specific parsing grammars, the parsing grammars for each culture are applied to each input record in the order the cultures are listed in the Open Parser stage. However, if you want to apply a specific culture's parsing grammar to a record, you can add a field named `CultureCode`. The field must contain one of the supported culture codes listed in the following table.

Culture Codes

Culture codes consist of a two-letter lowercase language code and a two-letter uppercase country or region code. For example, "es-MX" for Spanish (Mexico) and "en-US" for English (United States). In cases where a two-letter language code is not available, a three-letter code is used, for example "uz-Cyrl-UZ" for Uzbek (Uzbekistan, Cyrillic). A language is specified by only the two-digit lowercase language code. For example, "fr" specifies the neutral culture for French, and "de" specifies the neutral culture for German.

Note: There are two culture names that follow a different pattern. The cultures "zh-Hans" (Simplified Chinese) and "zh-Hant" (Traditional Chinese) are neutral cultures. The culture names represent the current standard and should be used unless you have a reason for using the older names "zh-CHS" and "zh-CHT".

The following table shows the supported culture codes.

Language (Culture/Region)	Culture Code
Global Culture	Global Culture
Afrikaans	af
Afrikaans (South Africa)	af-ZA
Albanian	sq

Language (Culture/Region)	Culture Code
Albanian (Albania)	sq-AL
Arabic	ar
Arabic (Algeria)	ar-DZ
Arabic (Bahrain)	ar-BH
Arabic (Egypt)	ar-EG
Arabic (Iraq)	ar-IQ
Arabic (Jordan)	ar-JO
Arabic (Kuwait)	ar-KW
Arabic (Lebanon)	ar-LB
Arabic (Libya)	ar-LY
Arabic (Morocco)	ar-MA
Arabic (Oman)	ar-OM
Arabic (Qatar)	ar-QA
Arabic (Saudi Arabia)	ar-SA

Language (Culture/Region)	Culture Code
Arabic (Syria)	ar-SY
Arabic (Tunisia)	ar-TN
Arabic (U.A.E.)	ar-AE
Arabic (Yemen)	ar-YE
Armenian	hy
Armenian (Armenia)	hy-AM
Azeri	az
Azeri (Azerbaijan, Cyrillic)	az-Cyrl-AZ
Azeri (Azerbaijan, Latin)	az-Latn-AZ
Basque	eu
Basque (Basque)	eu-ES
Belarusian	be
Belarusian (Belarus)	be-BY
Bulgarian	bg

Language (Culture/Region)	Culture Code
Bulgarian (Bulgaria)	bg-BG
Catalan	ca
Catalan (Catalan)	ca-ES
Chinese	zh
Chinese (Hong Kong SAR, PRC)	zh-HK
Chinese (Macao SAR)	zh-MO
Chinese (PRC)	zh-CN
Chinese (Simplified)	zh-Hans
Chinese (Singapore)	zh-SG
Chinese (Taiwan)	zh-TW
Chinese (Traditional)	zh-Hant
Croatian	hr
Croatian (Croatia)	hr-HR
Czech	cs

Language (Culture/Region)	Culture Code
Czech (Czech Republic)	cs-CZ
Danish	da
Danish (Denmark)	da-DK
Divehi	dv
Divehi (Maldives)	dv-MV
Dutch	nl
Dutch (Belgium)	nl-BE
Dutch (Netherlands)	nl-NL
English	en
English (Australia)	en-AU
English (Belize)	en-BZ
English (Canada)	en-CA
English (Caribbean)	en-029
English (Ireland)	en-IE

Language (Culture/Region)	Culture Code
English (Jamaica)	en-JM
English (New Zealand)	en-NZ
English (Philippines)	en-PH
English (South Africa)	en-ZA
English (Trinidad and Tobago)	en-TT
English (United Kingdom)	en-GB
English (United States)	en-US
English (Zimbabwe)	en-ZW
Estonian	et
Estonian (Estonia)	et-EE
Faroese	fo
Faroese (Faroe Islands)	fo-FO
Farsi	fa
Farsi (Iran)	fa-IR

Language (Culture/Region)	Culture Code
Finnish	fi
Finnish (Finland)	fi-FI
French	fr
French (Belgium)	fr-BE
French (Canada)	fr-CA
French (France)	fr-FR
French (Luxembourg)	fr-LU
French (Monaco)	fr-MC
French (Switzerland)	fr-CH
Galician	gl
Galician (Spain)	gl-ES
Georgian	ka
Georgian (Georgia)	ka-GE
German	de

Language (Culture/Region)	Culture Code
German (Austria)	de-AT
German (Germany)	de-DE
German (Liechtenstein)	de-LI
German (Luxembourg)	de-LU
German (Switzerland)	de-CH
Greek	el
Greek (Greece)	el-GR
Gujarati	gu
Gujarati (India)	gu-IN
Hebrew	he
Hebrew (Israel)	he-IL
Hindi	hi
Hindi (India)	hi-IN
Hungarian	hu

Language (Culture/Region)	Culture Code
Hungarian (Hungary)	hu-HU
Icelandic	is
Icelandic (Iceland)	is-IS
Indonesian	id
Indonesian (Indonesia)	id-ID
Italian	it
Italian (Italy)	it-IT
Italian (Switzerland)	it-CH
Japanese	ja
Japanese (Japan)	ja-JP
Kannada	kn
Kannada (India)	kn-IN
Kazakh	kk
Kazakh (Kazakhstan)	kk-KZ

Language (Culture/Region)	Culture Code
Konkani	kok
Konkani (India)	kok-IN
Korean	ko
Korean (Korea)	ko-KR
Kyrgyz	ky
Kyrgyz (Kyrgyzstan)	ky-KG
Latvian	lv
Latvian (Latvia)	lv-LV
Lithuanian	lt
Lithuanian (Lithuania)	lt-LT
Macedonian	mk
Macedonian (Macedonia, FYROM)	mk-MK
Malay	ms
Malay (Brunei Darussalam)	ms-BN

Language (Culture/Region)	Culture Code
Malay (Malaysia)	ms-MY
Marathi	mr
Marathi (India)	mr-IN
Mongolian	mn
Mongolian (Mongolia)	mn-MN
Norwegian	no
Norwegian (Bokmål, Norway)	nb-NO
Norwegian (Nynorsk, Norway)	nn-NO
Polish	pl
Polish (Poland)	pl-PL
Portuguese	pt
Portuguese (Brazil)	pt-BR
Portuguese (Portugal)	pt-PT
Punjabi	pa

Language (Culture/Region)	Culture Code
Punjabi (India)	pa-IN
Romanian	ro
Romanian (Romania)	ro-RO
Russian	ru
Russian (Russia)	ru-RU
Sanskrit	sa
Sanskrit (India)	sa-IN
Serbian	sr
Serbian (Serbia, Cyrillic)	sr-Cyrl-CS
Serbian (Serbia, Latin)	sr-Latn-CS
Slovak	sk
Slovak (Slovakia)	sk-SK
Slovenian	sl
Slovenian (Slovenia)	sl-SI

Language (Culture/Region)	Culture Code
Spanish	es
Spanish (Argentina)	es-AR
Spanish (Bolivia)	es-BO
Spanish (Chile)	es-CL
Spanish (Colombia)	es-CO
Spanish (Costa Rica)	es-CR
Spanish (Dominican Republic)	es-DO
Spanish (Ecuador)	es-EC
Spanish (El Salvador)	es-SV
Spanish (Guatemala)	es-GT
Spanish (Honduras)	es-HN
Spanish (Mexico)	es-MX
Spanish (Nicaragua)	es-NI
Spanish (Panama)	es-PA

Language (Culture/Region)	Culture Code
Spanish (Paraguay)	es-PY
Spanish (Peru)	es-PE
Spanish (Puerto Rico)	es-PR
Spanish (Spain)	es-ES
Spanish (Spain, Traditional Sort)	es-ES_tradnl
Spanish (Uruguay)	es-UY
Spanish (Venezuela)	es-VE
Swahili	sw
Swahili (Kenya)	sw-KE
Swedish	sv
Swedish (Finland)	sv-FI
Swedish (Sweden)	sv-SE
Syriac	syr
Syriac (Syria)	syr-SY

Language (Culture/Region)	Culture Code
Tamil	ta
Tamil (India)	ta-IN
Tatar	tt
Tatar (Russia)	tt-RU
Telugu	te
Telugu (India)	te-IN
Thai	th
Thai (Thailand)	th-TH
Turkish	tr
Turkish (Turkey)	tr-TR
Ukrainian	uk
Ukrainian (Ukraine)	uk-UA
Urdu	ur
Urdu (Pakistan)	ur-PK

Language (Culture/Region)	Culture Code
Uzbek	uz
Uzbek (Uzbekistan, Cyrillic)	uz-Cyrl-UZ
Uzbek (Uzbekistan, Latin)	uz-Latn-UZ
Vietnamese	vi
Vietnamese (Vietnam)	vi-VN

Grammars

A valid parsing grammar contains:

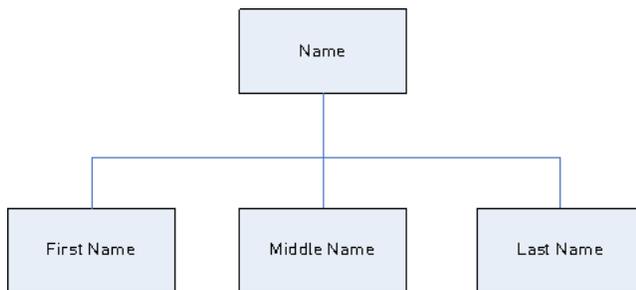
- A root variable that defines the sequence of tokens, or domain pattern, as rule variables.
- Rule variables that define the valid set of characters and the sequence in which those characters can occur in order to be considered a member of a domain pattern. For more information, see [Rule Section Commands](#) on page 30.
- The input field to parse. Input field designates the field to parse in the source data records.
- The output fields for the resulting parsed data. Output fields define where to store each resulting token that is parsed.

A valid parsing grammar also contains other optional commands for:

- Characters used to tokenize the input data that you are parsing. Tokenizing characters are characters, like space and hyphen, that determine the start and end of a token. The default tokenization character is a space. Tokenizing characters are the primary way that a sequence of characters is broken down into a set of tokens. You can set the tokenize command to NONE to stop the field from being tokenized. When tokenize is set to None, the grammar rules must include any spaces within its rule definition.
- Casing sensitivity options for tokens in the input data.
- Join character for delimiting matching tokens.
- Matching tokens in tables
- Matching compound tokens in tables
- Defining RegEx tags

- Literal strings in quotes
- Expression Quantifiers (optional). For more information about expression quantifiers, see [Rule Section Commands](#) on page 30 and [Expression Quantifiers: Greedy, Reluctant, and Possessive Behavior](#).
- Other miscellaneous indicators for grouping, commenting, and assignment (optional). For more information about grouped expressions, see [Grouping Operator\(\)](#).

The rule variables in your parsing grammar form a layered tree structure of the sequence of characters or tokens in a domain pattern. For example, you can create a parsing grammar that defines a domain pattern based on name input data that contains the tokens <FirstName>, <MiddleName>, and <LastName>.



Using the input data:

```
Joseph Arnold Cowers
```

You can represent that data string as three tokens in a domain pattern:

```
<root> = <FirstName><MiddleName><LastName>;
```

The rule variables for this domain pattern are:

```
<FirstName> = <given>;
<MiddleName> = <given>;
<LastName> = @Table("Family Names");
<given> = @Regex("[A-Za-z]+");
```

Based on this simple grammar example, Open Parser tokenizes on spaces and interprets the token `Joseph` as a first name because the characters in the first token match the `[A-Za-z]+` definition and the token is in the defined sequence. Optionally, any expression may be followed by another expression.

Example

```
<variable> = "some leading string" <variable2>;
<variable2> = @Table ("given") @Regex("[0-9]+");
```

A grammar rule is a grammatical statement wherein a variable is equal to one or more expressions. Each grammar rule follows the form:

```
<rule> = expression [| expression...];
```

Grammar rules must follow these rules:

- `<root>` is a special variable name and is the first rule executed in the grammar because it defines the domain pattern. `<root>` may not be referenced by any other rule in the grammar.
- A `<rule>` variable may not refer to itself directly or indirectly. When rule A refers to rule B, which refers to rule C, which refers to rule A, a circular reference is created. Circular references are not permitted.
- A `<rule>` variable is equal to one or more expressions.
- Each `expression` is separated by an OR, which is indicated using the pipe character "|".
- Expressions are examined one at a time. The first `expression` to match is selected. No further expressions are examined.
- The variable name may be composed of alphabetic, numeric, underscore (`_`) and hyphen (`-`). The name of the variable may start with any valid character. If the specified output field name does not conform to this form, use the alias feature to map the variable name to the output field.

An expression may be any of the following types:

- Another variable
- A string consisting of one or more characters in single or double quotes. For example:

```
"McDonald" 'McDonald' "O'Hara" 'O\'Hara' 'D"har' "D\"har"
```

- Table
- CompoundTable
- RegEx commands

Command Metacharacters

Open Parser supports the standard set of Java RegEx character class metacharacters in the `%Tokenize` and `@RegEx` commands. A metacharacter is a character that carries special meaning in pattern matching. The supported metacharacters are:

```
([{\^-$| ]})?*.+
```

There are two ways to force a metacharacter to be treated as an ordinary character:

- Precede the metacharacter with a backslash
- Enclose it within `\Q` (which starts the quote) and `\E` (which ends it).

`%Tokenize` follows the rule for Java Regular Expressions character classes—not Java Regular Expressions as a whole.

In general, the reserved characters for a character set are:

- '[' and ']' indicate another set.
- '-' is a metacharacter if in between two other characters.
- '^' is a metacharacter if it is the first character in a set.
- '&&' are metacharacters if they are between two other characters.

- `\` means next that the character is a literal.

If you have any doubt whether a character will be treated as a metacharacter and you want the character to be treated as a literal, escape that character using the backlash.

Header Section Commands

This section describes the header section commands. Some commands are optional. If a command is optional, the default value or behavior is listed.

- **Tokenize** (optional)
- **Tokenize (None)**
- **InputField** (required if Input Fields is not used)
- **InputFields** (required if Input Field is not used)
- **OutputFields** (required)
- **IgnoreCase** (optional)
- **Join** (optional)

Rule Section Commands

The rule section commands are:

- **RegEx**
- **Table**
- **CompoundTable**
- **Token**
- **Scoring**
- **RuleID**
- **<root> Variable**
- **rule|rule**
- **Grouping Operator ()**
- **Min/Max Occurrences Operator {min,max}**
- **Exact Occurrences Operator {exact}**
- **Assignment Operator (=)**
- **OR Operator (|)**
- **End-of-Rule Operator (;)**
- **Commenting Operator (!)**
- **Zero or One Occurrences Quantifier (?)**
- **Zero or More Occurrences Quantifier (*)**
- **One or More Occurrences Quantifier (+)**
- **Expression Quantifiers: Greedy, Reluctant, and Possessive Behavior**

Cultures

A culture is the primary concept for organizing culture-specific parsing grammars. You can use cultures to create different parsing rules for different cultures and languages. Culture follows a hierarchy:

- **Global Culture:** The global culture is culture-independent and language agnostic. Use global culture to create parsing grammar rules that span all cultures and languages.
- **Language:** A language is associated with a language, but not with a specific culture/region. For example, English.
- **Culture/Region:** A culture/region is associated with a language and a country or region. For example, English in the United Kingdom, or English in the United States.

In the culture hierarchy, the parent of a culture/region is a language and the parent of a language is the global culture.

Culture/regions inherit the properties of the parent language. Languages inherit the properties of the global culture. As such, you can define parsing grammars in a language for use in multiple countries that share that language. Then, you can override the language grammar rules with specialized parsing grammars for a particular country or region that shares the same language as the base language culture, but has specific addressing, naming, or other country or regional differences.

You can also use culture inheritance to parse incoming records that have an assigned culture code, but no defined grammar rule for that culture code. In this case, Open Parser looks for a language code that has an assigned grammar rule. If it does not exist, Open Parser looks for an assigned grammar rule in the global culture.

The Domain Editor uses a combination of a language code and a culture code to represent language and culture/region, respectively.

Defining a Culture's Grammar Rules

You can use a culture's grammar rules to substitute a portion of a the global culture's parsing grammar with strings, commands, or expressions specific to the culture and/or language. By defining a grammar rule, you can customize portions of the global culture parsing grammar based on the record's culture and/or language. This is useful if you do not want to create an entirely separate parsing grammar for each culture and instead use the global culture's grammar, customizing only specific portions of the global culture grammar for each culture.

This topic describes how to create a grammar rule for a culture.

1. In Enterprise Designer, go to **Tools > Open Parser Domain Editor**.
2. Click the **Cultures** tab.

For a complete list of supported cultures, see [Assigning a Parsing Culture to a Record](#) on page 12.

3. Select the culture to which you want to add a grammar rule then click **Properties**.
4. Click the **Grammar Rules** tab. The information displayed includes the grammar rule names defined for the selected culture, the associated source culture, the defined value of the grammar rule, and the description.
5. Click **Add**.
6. Type a name for the grammar rule in the **Name** field.
7. Type a description of the grammar rule in the **Description** field.
8. Type the grammar rule in the **Value** field.

The grammar rule can be any valid variable, string, command, or grouped expression. For more information, see [Grammars](#) on page 27.

9. Select **Enable word wrap** to display the value in the text box without scrolling.
10. Click **OK**.

The grammar rule value that you typed is validated. If the value contains grammar syntax errors, a message displays a description of the errors encountered, the line and column where the error occurs, and the command, grammar rule, or RegEx tag where the error occurs.

Example Grammar Rule

You have a grammar that parses Western names. The structure of the pattern maybe the same for all cultures (<FirstName><MiddleName><LastName>) and many of the rules might match the same pattern or table. However, you also have culture-specific tables for last names, and you want to use the appropriate table based on the record's culture code.

To accomplish this, you could define a grammar rule for each culture that replaces the <LastName> element in the global culture with a reference to the culture-specific table. For example, if you have a table of Dutch last names, you would create a grammar rule for the Dutch (nl) culture as follows:

Name: LastName
Description: Dutch last names
Value: @Table("Dutch Last Names");

Defining Culture RegEx Tags

This topic describes how to define culture RegEx tags when defining a culture-specific parsing grammar.

1. In Enterprise Designer, go to **Tools > Open Parser Domain Editor**.

2. Click the **Cultures** tab. The **Cultures** tab displays a list of supported cultures. For a complete list of supported cultures, see [Assigning a Parsing Culture to a Record](#) on page 12.
3. Select a culture from the list and then click **Properties**. The **Culture Properties** dialog box displays.
4. Click the **RegEx Tags** tab. The information displayed includes the RegEx tag names defined for the selected culture and the associated source culture, the value of the RegEx tag, and the description.
5. Click **Add** or **Modify**.
6. Type a name for the RegEx tag in the **Name** text box.

If you type a name that already exists in the selected culture, a warning icon flashes. Type a different name or close the dialog box, delete the existing RegEx tag, and then click **Add** again.

7. Type a description of the RegEx tag in the **Description** text box.
8. Type a value for the RegEx tag in the **Value** text box.

The value can be any valid regular expression but cannot match an empty string.

Domain Editor includes several predefined RegEx tags that you can use to define culture properties. You can also use these RegEx tags for defining tokenization characters in your parsing grammar.

You can modify the predefined RegEx tags or copy them and create your own variants. You can also use override properties to create specialized RegEx tags for specific languages.

- **Letter**: Any letter from any language. This RegEx tag includes overrides for several languages due to differences in scripts used, for example, cyrillic scripts, asian-language scripts, and Thai script.
- **Lower**: A lowercase letter that has an uppercase variant.
- **Number**: Any numeric character in any script.
- **Punctuation**: Any punctuation character.
- **Upper**: An uppercase letter that has a lowercase variant.
- **Whitespace**: Any whitespace or invisible separator.

9. Click **OK**.

Importing and Exporting Cultures

In addition to creating cultures, you can also import cultures you've created elsewhere and export cultures you create in the Domain Editor.

1. In Enterprise Designer, go to **Tools > Open Parser Domain Editor**.
2. Click the **Cultures** tab.
3. Click **Import** or **Export**.
4. Do one of the following:

- If you are importing a culture, navigate to and select a culture. Click **Open**. The imported culture appears in the Domain Editor.
- If you are exporting a culture, navigate to and select the location where you would like to save the exported culture. Click **Save**. The exported culture is saved and the Domain Editor returns.

Domains

Adding a Domain

A domain represents a type of data such as name, address, and phone number data. It consists of a pattern that represents a sequence of one or more tokens in your input data that you commonly need to parse and that you associate with one or more cultures.

This topic describes how to add a domain in Domain Editor when defining a culture-specific parsing grammar.

1. In Enterprise Designer, go to **Tools > Open Parser Domain Editor**.
2. Click the **Domains** tab.
3. Click **Add**.
4. Type a domain name in the **Name** field.
5. Type a description of the domain name in the **Description** field.
6. If you want to create a new, empty domain, click **OK**. If you want to create a new domain based on another domain, do the following:
 - a) Select **Use another domain as a template** if you want to create a new domain based on another domain.
 - b) Select a domain from the list. When you click **OK** in the next step, the new domain will be created. The new domain will contain all of the culture-specific parsing grammars defined in the domain template that you selected.
 - c) Click **OK**.

Modifying a Domain

A domain represents a type of data such as name, address, and phone number data. It consists of a pattern that represents a sequence of one or more tokens in your input data that you commonly need to parse and that you associate with one or more cultures.

This topic describes how to modify a domain.

1. In Enterprise Designer, go to **Tools > Open Parser Domain Editor**.
2. Click the **Domains** tab.
3. Select a domain in the list and then click **Modify**. The **Modify Domain** dialog box displays.

4. Change the description information.
5. If you only want to modify the description of the domain, click **OK**. If you have made updates to the template domain and now want to add those changes to the domain you are modifying, then continue to the next step.
6. Select **Use another domain as a template** to inherit changes made to the domain template.
7. Select a domain pattern template from the list. When you click **OK** in the next step, the domain pattern will be modified. The modified domain pattern will contain all of the culture-specific parsing grammars defined in the domain pattern template that you selected. Any parsing grammar in the selected domain pattern will be overwritten with the parsing grammar from the domain pattern template.
8. Click **OK**.

To see how this works, do the following:

1. Create a domain pattern named **NameParsing** and define parsing grammars for Global Culture, en, and en-US.
2. Create a domain pattern named **NameParsing2** and use **NameParsing** as a domain pattern template. **NameParsing2** is created as an exact copy and contains parsing grammars for Global Culture, en, and en-US.
3. Modify the culture-specific parsing grammars for **NameParsing** by changing some of the grammar rules in the Global Culture grammar and add en-CA as a new culture.
4. Select **NameParsing2** on the Domains tab, click Modify, and again use **NameParsing** as the domain pattern template.

The results will be:

- The Global Culture parsing grammar will be updated (overwriting your changes if any have been made).
- The cultures en and en-US will remain the same (unless they have been modified in the target domain, in which case they would then revert back to the Name Parsing version).
- A culture-specific grammar for en-CA will be added.

Removing a Domain

A domain represents a type of data such as name, address, and phone number data. It consists of a pattern that represents a sequence of one or more tokens in your input data that you commonly need to parse and that you associate with one or more cultures.

This topic describes how to remove a domain.

1. In Enterprise Designer, go to **Tools > Open Parser Domain Editor**.
2. Click the **Domains** tab.
3. Select a domain in the list.

4. Click **Remove**.

If the domain is associated with one or more culture-specific parsing grammars, a message displays asking you to confirm that you want to remove the domain. If no culture-specific parsing grammars are associated with this domain, a message displays confirming that you want to remove the selected domain.

5. Click **Yes**. The domain and any culture-specific parsing grammars associated with this domain are removed.

Importing and Exporting Domains

In addition to creating domains, you can also import domains you've created elsewhere and export domains you create in the Domain Editor.

1. Click the **Domains** tab. The **Domains** tab displays.
2. Click **Import** or **Export**.
3. Do one of the following:
 - If you are importing a domain, navigate to and select a domain name. Click **Open**. The imported domain appears in the Domain Editor.
 - If you are exporting a domain, navigate to and select the location where you would like to save the exported domain. Click **Save**. The exported domain is saved and the Domain Editor returns.

Analyzing Parsing Results

Tracing Final Parsing Results

The Open Parser Trace Details feature displays a graphical view of how the input field was parsed, token-by-token, into the output field values. Trace displays matching results, non-matching results, and interim results.

Final Parsing Results shows the parsing grammar tree and the resulting output. Use this view when you want to see only the results of the matching process. This is the default view.

1. In Enterprise Designer, open the dataflow that contains the Open Parser stage whose parsing results you want to trace.
2. Double-click the Open Parser stage on the canvas.
3. Click the **Preview** tab.
4. Enter sample data that you want to parse then click the **Preview** button.

- In the Trace column, click the **Click here...** link to display the trace diagram.

The tree view of the parsing grammar shows one or more the following elements, depending on the selected options:

- The <root> variable. The top node in the tree is the <root> variable.
- The expressions defined in the <root> variable. The second-level nodes are the expressions defined in the <root> variable. The <root> expressions also define the names of the output fields.
- The variable definitions of the second-level nodes. The third-level nodes and each level below it are the definitions of each of the <root> expressions. Expression definitions can be other variables, aliases, or rule definitions.
- The values and tokens that are output. The bottom node in the tree shows the values assigned to each sequential token in the parsing grammar.
- The parser score for relevant elements of the parsing grammar. Parser scores are determined from the bottom of a root expression to the top. For example, if an expression pattern has a weight of 80 and an ancestor rule has a weight of 75, the final score for the ancestor expression is the product of the child scores and the ancestor scores, which in this example would be 60 percent.
- The space character displays in the **Input data** text box as a non-breaking space character (upward facing bracket) so that you can better see space characters. Delimiters not used as tokens are displayed as gray.

- In the **Information** field, select **Final parsing results**.

Note: To step through the parsing events, see [Stepping Through Parsing Events](#) on page 38.

- In the **Level of detail** list, select one of the options.

- **Hide expressions without results.** Shows those branches that lead to a matching or non-matching result. Any root expression branch that does not lead to a match is shown as an ellipsis. If you want to look at a branch that does not lead to a match, double-click on the ellipsis.
- **Hide root expressions without results.** Shows all branches of the root expressions containing match or non-matching results. Any other root expressions are not displayed.
- **Show all roots.** Shows every root expression. If a root has no matching result, the display is collapsed for that root expression using the ellipsis symbol.
- **Show all expressions.** Shows the root expressions and all branches. The root expressions are no longer displayed as an ellipsis; instead, the rules for each expression in the branch are shown.

If you have a level-of-detail view selected that hides expressions without results and you select a root expression that is not currently displayed, Trace Details changes the level-of-detail selection to a list item that shows the minimum number of root expressions, while still displaying the root expression.

8. Click **Show scores** to display parser scores for root expressions, variable expressions, and the resulting matches and non-matches.
9. In the **Zoom** field, select the size of the tree view.
10. In the **Root clause** field, select one of the options to show that branch of the root expression tree.

When you click an expression branch in the trace diagram, the **Root clause** list updates to display the selected clause. Double-click an ellipsis to display a collapsed expression.

11. Click **OK** when you are done. The level of detail, show scores, and zoom control settings are saved when you click **OK**.

Stepping Through Parsing Events

The Open Parser Trace Details view allows you to view a diagram of event-by-event steps in the matching process. Use this view when you are troubleshooting the matching process and want to see how each token is evaluated, the parsing grammar tokenization, and the token-by-token matching results.

1. In Enterprise Designer, open the dataflow that contains the Open Parser stage whose parsing results you want to trace.
2. Double-click the Open Parser stage on the canvas.
3. Click the **Preview** tab.
4. Enter sample data that you want to parse then click the **Preview** button.
5. In the Trace column, click the **Click here...** link to display the trace diagram.

The tree view of the parsing grammar shows one or more the following elements, depending on the selected options:

- The <root> variable. The top node in the tree is the <root> variable.
- The expressions defined in the <root> variable. The second-level nodes are the expressions defined in the <root> variable. The <root> expressions also define the names of the output fields.
- The variable definitions of the second-level nodes. The third-level nodes and each level below it are the definitions of each of the <root> expressions. Expression definitions can be other variables, aliases, or rule definitions.
- The values and tokens that are output. The bottom node in the tree shows the values assigned to each sequential token in the parsing grammar.
- The parser score for relevant elements of the parsing grammar. Parser scores are determined from the bottom of a root expression to the top. For example, if an expression pattern has a weight of 80 and an ancestor rule has a weight of 75, the final score for the ancestor expression is the product of the child scores and the ancestor scores, which in this example would be 60 percent.

- The space character displays in the **Input data** text box as a non-breaking space character (upward facing bracket) so that you can better see space characters. Delimiters not used as tokens are displayed as gray.
- Matches and non-matches are color coded in the trace diagram:
 - Green boxes indicate matches that are part of the final successful result.
 - Red boxes indicate non-matches.
 - Yellow boxes indicate interim matches that will eventually be rolled back as the events are stepped through. Interim matches display only in Step Through Parsing Events.
 - Gray boxes indicate interim matches that have been rolled back to free up that token for another expression. Interim matches display only in Step Through Parsing Events.
 - In the **Information** list, select **Step through parsing events**.
 - In the **Level of detail** list, select one of the options.

- **Hide expressions without results.** Shows those branches that lead to a matching or non-matching result. Any root expression branch that does not lead to a match is shown as an ellipsis. If you want to look at a branch that does not lead to a match, double-click on the ellipsis.
- **Hide root expressions without results.** Shows all branches of the root expressions containing match or non-matching results. Any other root expressions are not displayed.
- **Show all roots.** Shows every root expression. If a root has no matching result, the display is collapsed for that root expression using the ellipsis symbol.
- **Show all expressions.** Shows the root expressions and all branches. The root expressions are no longer displayed as an ellipsis; instead, the rules for each expression in the branch are shown.

If you have a level-of-detail view selected that hides expressions without results and you select a root expression that is not currently displayed, Trace Details changes the level-of-detail selection to a list item that shows the minimum number of root expressions, while still displaying the root expression.

- Click **Show scores** to display parser scores for root expressions, variable expressions, and the resulting matches and non-matches.
- In the **Zoom** field, select the size of the tree view.
- In the **Root clause** field, select one of the options to show that branch of the root expression tree.

When you click an expression branch in the trace diagram, the **Root clause** list updates to display the selected clause. Double-click an ellipsis to display a collapsed expression.

- The **Automatically step to selected node** check box is selected by default. When this is selected and you click the **Play** button, the events execute from the beginning and stop on the first event that occurs with the selected node or any of its children. To play all events without stopping, clear this check box before clicking the **Play** button.
- In the **Play delay (seconds)** field, specify a delay to control the speed of the play rate.

14. Click the **Play** button to start executing the parsing events.
15. Click **OK** when you are done.

Parsing Personal Names

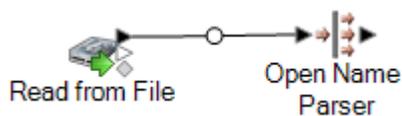
If you have name data that is all in one field, you may want to parse the name into separate fields for each part of the name, such as first name, last name, title of respect, and so on. These parsed name elements can then be used by other automated operations such as name matching, name standardization, or multirecord name consolidation.

1. If you have not already done so, load the following tables onto the Spectrum Technology Platform server:
 - Open Parser Base
 - Open Parser Enhanced Names

Use the Data Normalization database load utility to load these tables. For instructions on loading tables, see the *Installation Guide*.

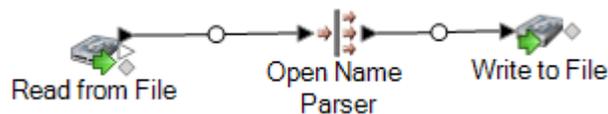
2. In Enterprise Designer, create a new dataflow.
3. Drag a source stage onto the canvas.
4. Double-click the source stage and configure it. See the *Dataflow Designer's Guide* for instructions on configuring source stages.
5. Drag an Open Name Parser stage onto the canvas and connect it to the source stage.

For example, if you are using a Read from File stage, your dataflow would look like this:



6. Drag a sink stage onto the canvas and connect Open Name Parser to it.

For example, if you are using a Write to File sink, your dataflow might look like this:



7. Double-click the sink stage and configure it. See the *Dataflow Designer's Guide* for instructions on configuring source stages.

You have created a dataflow that can parse personal names into component parts, placing each part of the name in its own field.

Dataflow Templates for Parsing

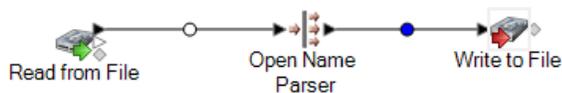
Parsing English Names

This dataflow template demonstrates how to take personal name data (for example "John P. Smith"), parse it into first name, middle name, and last name parts, and add gender data.

Business Scenario

You work for an insurance company that wants to send out personalized quotes based on gender to prospective customers. Your input data include name data as full names and you want to parse the name data into First, Middle, and Last name fields. You also want to determine the gender of the individuals in your input data.

The following dataflow provides a solution to the business scenario:



This dataflow template is available in Enterprise Designer. Go to **File > New > Dataflow > From template** and select **Parse Personal Name**.

This dataflow requires the following:

- The Universal Name
- The Open Parser base tables
- The Open Parser enhanced names tables

In this dataflow, data is read from a file and processed through the Open Name Parser stage. Open Name Parser is a part of Universal Naming. For each name, the dataflow does the following:

Read from File

This stage identifies the file name, location, and layout of the file that contains the names you want to parse. The file contains both male and female names.

Open Name Parser

Open Name Parser examines name fields and compares them to name data stored in the Spectrum Technology Platform name database files. Based on the comparison, it parses the name data into First, Middle, and Last name fields.

Write to File

The template contains one Write to File stage. In addition to the input fields, the output file contains the FirstName, MiddleName, LastName, EntityType, GenderCode, and GenderDeterminationSource fields.

Parsing Arabic Names

This template demonstrates how to parse westernized Arabic names into component parts. The parsing rule separates each token in the **Name** field and copies each token to five fields: **Kunya**, **Ism**, **Laqab**, **Nasab**, **Nisba**. These output fields represent the five parts of an Arabic name and are described in the business scenario.

Business Scenario

You work for a bank that wants to better understand the Arabic naming system in an effort to improve customer service with Arabic-speaking customers. You have had complaints from customers whose billing information does not list the customer's name accurately. In an effort to improve customer intimacy, the Marketing group you work in wants to better address Arabic-speaking customers through marketing campaigns and telephone support.

In order to understand the Arabic naming system, you search for and find these resources on the internet that explain the Arabic naming system:

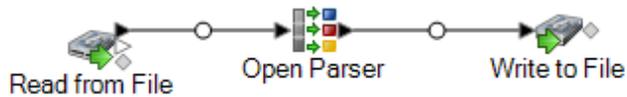
- en.wikipedia.org/wiki/Arabic_names
- heraldry.sca.org/laurel/names/arabic-naming2.htm

Arabic names are based on a naming system that includes these name parts: Ism, Kunya, Nasab, Laqab, and Nisba.

- The ism is the main name, or personal name, of an Arab person.
- Often, a kunya referring to the person's first-born son is used as a substitute for the ism.
- The nasab is a patronymic or series of patronymics. It indicates the person's heritage by the word ibn or bin, which means son, and bint, which means daughter.
- The laqab is intended as a description of the person. For example, al-Rashid means the righteous or the rightly-guided and al-Jamil means beautiful.
- The nisba describes a person's occupation, geographic home area, or descent (tribe, family, and so on). It will follow a family through several generations. The nisba, among the components of

the Arabic name, perhaps most closely resembles the Western surname. For example, al-Filistin means the Palestinian.

The following dataflow provides a solution to the business scenario:



This dataflow template is available in Enterprise Designer. Go to **File > New > Dataflow > From template** and select **ParseArabicNames**. This dataflow requires Data Normalization.

In this dataflow, data is read from a file and processed through the Open Parser stage. For each data row in the input file, this dataflow will do the following:

Read from File

This stage identifies the file name, location, and layout of the file that contains the names you want to parse. The file contains both male and female names.

Open Parser

This stage defines whether to use a culture-specific domain grammar created in the Domain Editor or to define a domain-independent grammar. A culture-specific parsing grammar that you create in the Domain Editor is a validated parsing grammar that is associated with a culture and a domain. A domain-independent parsing grammar that you create in Open Parser is a validated parsing grammar that is not associated with a culture and domain.

In this template, the parsing grammar is defined as a domain-independent grammar.

The Open Parser stage contains a parsing grammar that defines the following commands and expressions:

- `%Tokenize` is set to the space character (`\s`). This means that Open Parser will use the space character to separate the input field into tokens. For example, Abu Mohammed al-Rahim ibn Salamah contains five tokens: Abu, Mohammed, al-Rahim, ibn and Salamah.
- `%InputField` is set to parse input data from the **Name** field.
- `%OutputFields` is set to copy parsed data into five fields: **Kunya**, **Ism**, **Laqab**, **Nasab**, and **Nisba**.
- The `<root>` expression defines the pattern for Arabic names:
- Zero or one occurrence of **Kunya**
- Exactly one or two occurrences of **Ism**
- Zero or one occurrence of **Laqab**
- Zero or one occurrence of **Nasab**
- Zero or more occurrences of **Nisba**

The rule variables that define the domain must use the same names as the output fields defined in the required `OutputFields` command.

The parsing grammar uses a combination of regular expressions and expression quantifiers to build a pattern for Arabic names. The parsing grammar uses these special characters:

- The "?" character means that a regular expression can occur zero or one time.
- The "*" character means that a regular expression can occur zero or more times
- The ";" character means end of a rule.

Use the **Commands** tab to explore the meaning of the other special symbols you can use in parsing grammars by hovering the mouse over the description.

By default, quantifiers are greedy. Greedy means that the expression accepts as many tokens as possible, while still permitting a successful match. You can override this behavior by appending a '?' for reluctant matching or '+' for possessive matching. Reluctant matching means that the expression accepts as few tokens as possible, while still permitting a successful match. Possessive matching means that the expression accepts as many tokens as possible, even if doing so prevents a match.

To test the parsing grammar, click the Preview tab. Type the names shown below in the **Name** field and then click **Preview**.

Name	Kunya	Ism	Laqab	Nasab	Nisba
Abu Karim Muhammad al-Jamil ibn Nidal ibn Abdulaziz al-Filistini	Abu Karim	Muhammad	al-Jamil	ibn Nidal ibn Abdulaziz	al-Filistini
Layla bint Zuhayr ibn Yazid al-Nahdiyah		Layla		bint Zuhayr ibn Yazid	al-Nahdiyah
Yazid ibn Abi Hakim		Yazid		ibn Abi Hakim	
Abu Bishr al-Yaman ibn Abi al-Yaman al-Bandaniji	Abu Bishr	al-Yaman		ibn Abi	al-Yaman al-Bandaniji
Abu al-Tayyib 'Abd al-Rahim ibn Ahmad al-Harrani	Abu al-Tayyib	'Abd	al-Rahim	ibn Ahmad	al-Harrani
Ahmad ibn Sa'id al-Bahili		Ahmad		ibn Sa'id	al-Bahili
Abu al-Abbas Muhammad ibn Ya'qub ibn Yusuf al-Asamm al-Naysaburi	Abu al-Abbas	Muhammad		ibn Ya'qub ibn Yusuf	al-Asamm al-Naysaburi
Abu al-Qasim Mansur ibn al-Zabriqan ibn Salamah al-Namari	Abu al-Qasim	Mansur		ibn al-Zabriqan ibn Salamah	al-Namari
'Ubayd ibn Mu'awiyah ibn Zayd ibn Thabit ibn al-Dahhak		'Ubayd		ibn Mu'awiyah ibn Zayd ibn Thabit ibn al-Dahhak	
Umm Ja'far Zubaydah	Umm Ja'far	Zubaydah			

You can also type other valid and invalid names to see how the input data is parsed.

You can use the Trace feature to see a graphical representation of either the final parsing results or to step through the parsing events. Click the link in the **Trace** column to see the Trace Details for the data row.

Write to File

The template contains one Write to File stage. In addition to the input field, the output file contains the **Kunya, Ism, Laqab, Nasab, and Nisba** fields.

Parsing Chinese Names

This template demonstrates how to parse Chinese names into component parts. The parsing rule separates each token in the **Name** field and copies each token to two fields: **LastName** and **FirstName**.

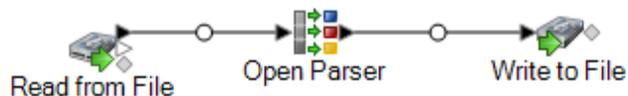
Business Scenario

You work for a financial service company that wants to explore if it is feasible to include the Chinese characters for its Chinese-speaking customers on various correspondence.

In order to understand the Chinese naming system, you search for and find this resource on the internet, which explains how Chinese names are formed:

en.wikipedia.org/wiki/Chinese_names

The following dataflow provides a solution to the business scenario:



This dataflow template is available in Enterprise Designer. Go to **File > New > Dataflow > From template** and select **ParseChineseNames**. This dataflow requires Data Normalization.

In this dataflow, data is read from a file and processed through the Open Parser stage. For each data row in the input file, this data flow will do the following:

Read from File

This stage identifies the file name, location, and layout of the file that contains the names you want to parse. The file contains both male and female names.

Open Parser

This stage defines whether to use a culture-specific domain grammar created in the Domain Editor or to define a domain-independent grammar. A culture-specific parsing grammar that you create in the Domain Editor is a validated parsing grammar that is associated with a culture and a domain. A domain-independent parsing grammar that you create in Open Parser is a validated parsing grammar that is not associated with a culture and domain.

In this template, the parsing grammar is defined as a domain-independent grammar.

The Open Parser stage contains a parsing grammar that defines the following commands and expressions:

- `%Tokenize` is set to `None`. When `Tokenize` is set to `None`, the parsing grammar rule must include any spaces or other token separators within its rule definition.
- `%InputField` is set to parse input data from the **Name** field.
- `%OutputFields` is set to copy parsed data into two fields: **LastName** and **FirstName**.

The `<root>` expression defines the pattern for Chinese names:

- One occurrence of **LastName**
- One to three occurrences of **FirstName**

The rule variables that define the domain must use the same names as the output fields defined in the required `OutputFields` command.

The `CJKCharacter` rule variable defines the character pattern for Chinese/ Japanese/Korean (CJK). The character pattern is defined so as to only use characters that are letters. The rule is:

```
<CJKCharacter> = @Regex("([\p{InCJKUnifiedIdeographs}&&\p{L}])");
```

- The regular expression `\p{InX}` is used to indicate a Unicode block for a certain culture, in which `X` is the culture. In this instance the culture is `CJKUnifiedIdeographs`.
- In regular expressions, a character class is a set of characters that you want to match. For example, `[aeiou]` is the character class containing only vowels. Character classes may appear within other character classes, and may be composed by the union operator (implicit) and the intersection operator (`&&`). The union operator denotes a class that contains every character that is in at least one of its operand classes. The intersection operator denotes a class that contains every character that overlaps the intersected Unicode blocks.
- The regular expression `\p{L}` is used to indicate the Unicode block that includes only letters.

To test the parsing grammar, click the Preview tab. Type the names shown below in the **Name** field and then click **Preview**.

Name	FirstName	LastName
王若琳	若琳	王
刘翀宏	翀宏	刘
许惠欣	惠欣	许
蔡依林	依林	蔡
水沫	沫	水
龚亚轩	亚轩	龚
郑元畅	元畅	郑
林依晨	依晨	林

You can also type other valid and invalid names to see how the input data is parsed.

You can use the Trace feature to see a graphical representation of either the final parsing results or to step through the parsing events. Click the link in the **Trace** column to see the Trace Details for the data row.

Write to File

The template contains one Write to File stage. In addition to the input field, the output file contains the **LastName**, and **FirstName** fields. Select a match results in the Match Results List and then click **Remove**.

Parsing Spanish and German Names

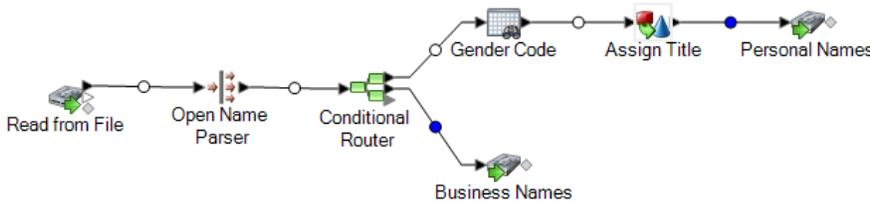
This template demonstrates how to parse mixed-culture names, such as Spanish and German names, into component parts. The parsing rule separates each token in the **Name** field and copies each token to the fields defined in the Personal and Business Names parsing grammar. For more information about this parsing grammar, select **Tools > Open Parser Domain Editor** and then select the **Personal and Business Names** domain and either the **German (de)** or **Spanish (es)** cultures.

This template also applies gender codes to personal names in using table data contained in Table Management. For more information about Table Management, select **Tools > Table Management**.

Business Scenario

You work for a pharmaceuticals company based in Brussels that has consolidated its Germany and Spain operations. Your company wants to implement a mixed-culture database containing name data and it is your job to analyze the variations in names between the two cultures.

The following dataflow provides a solution to the business scenario:



This dataflow template is available in Enterprise Designer. Go to **File > New > Dataflow > From template** and select **ParseSpanish&GermanNames**. This dataflow requires Data Normalization.

In this dataflow, data is read from a file and processed through the Open Parser stage. For each data row in the input file, this data flow will do the following:

Read from File

This stage identifies the file name, location, and layout of the file that contains the names you want to parse. The file contains both male and female names and includes CultureCode information for each name. The CultureCode information designates the input names as either German (de) or Spanish (es).

Open Name Parser

Open Name Parser examines name fields and compares them to name data stored in the Spectrum Technology Platform name database files. Based on the comparison, it parses the name data into First, Middle, and Last name fields.

Conditional Router

This stage routes the input so that personal names are routed to the Gender Codes stage and business names are routed to the Business Names stage.

Gender Code

Double-click this stage on the canvas and then click **Modify** to display the table lookup rule options.

The **Categorize** option uses the Source value as a key and copies the corresponding value from the table entry into the field selected in the Destination list. In this template, **Complete field** is selected and **Source** is set to use the **FirstName** field. Table Lookup treats the entire field as one string and flags the record if the string as a whole can be categorized.

The **Destination** is set to the **GenderCode** field and uses the lookup terms contained in the **Gender Codes** table to perform the categorization of male and female names. If a term in the input data is not found, Table Lookup assigns a value of **U**, which means unknown. To better understand how this works, select **Tools > Table Management** and select the Gender Codes table.

Write to File

The template contains two Write to File stages, one for personal names and one for business names. In addition to the input field, the personal names output file contains the **Name**, **TitleOfRespect**, **FirstName**, **MiddleName**, **LastName**, **PaternalLastName**, **MaternalLastName**, **MaturitySuffix**, **GenderCode**, **CultureUsed**, and **ParserScore** fields.

The business names output file contains the **Name**, **FirmName**, **FirmSuffix**, **CultureUsed**, and **ParserScore** fields.

Parsing E-mail Addresses

This template demonstrates how to parse e-mail addresses into component parts. The parsing rule separates each token in the **Email** field and copies each token to three fields: **Local-Part**, **DomainName**, and **DomainExtension**. **Local-Part** represents the domain name part of the e-mail address, **DomainName** represents the domain name of the e-mail address, and **DomainExtension** represents the domain extension of the e-mail address. For example, in `precisely.com`, "precisely" is the domain name and "com" is the domain extension.

The internet is a great source of public domain information that can aid you in your open parsing tasks. In this example, e-mail formatting information was obtained from various internet resources and was then imported into Table Management to create a table of domain values. The domain extension task that you will perform in this template activity demonstrates the usefulness of this method.

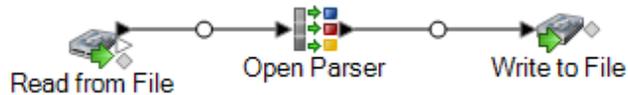
This template also demonstrates how to effectively use table data that you load into Table Management to perform table look-ups as part of your parsing tasks.

Business Scenario

You work for an insurance company that wants to do its first e-mail marketing campaign. Your database contains e-mail addresses of your customers and you have been asked to find a way to make sure that those e-mail addresses are in a valid SMTP format.

Before you create this dataflow, you will need to load a table of valid domain names extensions in Table Management so that you can look up domain name extensions as part of the validation process.

The following dataflow provides a solution to the business scenario:



This dataflow template is available in Enterprise Designer. Go to **File > New > Dataflow > From template** and select **ParseEmail**. This dataflow requires Data Normalization.

In this dataflow, data is read from a file and processed through the Open Parser stage. For each data row in the input file, this dataflow will do the following:

Create a Domain Extension Table

The first task is to create an Open Parser table in Table Management that you can use to check if the domain extensions in your e-mail addresses are valid.

1. From the **Tools** menu, select **Table Management**.
2. In the **Type** list, select **Open Parser**.
3. Click **New**.
4. In the **Add User Defined Table** dialog box, type `EmailDomains` in the **Table Name** field, make sure that **None** is selected in the **Copy from** list, and then click **OK**.
5. With **EmailDomains** displayed in the **Name** list, click **Import**.
6. In the **Import** dialog box, click **Browse** and locate the source file for the table. The default location is: `<drive>:\Program Files\Precisely\Spectrum\server\modules\coretemplates\data\Email_Domains.txt`. Table Management displays a preview of the terms contained in the import file.
7. Click **OK**. Table Management imports the source files and displays a list of internet domain extensions.
8. Click **Close**. The `EmailDomains` table is created. Now create the dataflow using the ParseEmail template.

Read from File

This stage identifies the file name, location, and layout of the file that contains the e-mail addresses you want to parse.

Open Parser

The Open Parser stage parsing grammar defines the following commands and expressions:

- %Tokenize is set to None. When Tokenize is set to None, the parsing grammar rule must include any spaces or other token separators within its rule definition.
- %InputField is set to parse input data from the **Email_Address** field.
- %OutputFields is set to copy parsed data into three fields: **Local-Part**, **DomainName**, and **DomainExtension**.
- The root expression defines the pattern of tokens being parsed:

```
<root> = <Local-Part>@"<DomainName>".<DomainExtension>;
```

The rule variables that define the domain must use the same names as the output fields defined in the required `OutputFields` command.

- The remainder of the parsing grammar defines each of the rule variables as expressions.

```
<Local-Part> = (<alphanum> ".")* <alphanum> | (<alphanum> "_")* <alphanum>
;
<DomainName> = (<alphanum> ".")? <alphanum>;
<DomainExtension> = @Table("EmailDomains")* "."? @Table("EmailDomains");
<alphanum>=@RegEx("[A-Za-z0-9]+");
```

The `<Local-Part>` variable is defined as a string of text that contains the `<alphanum>` variable, the period character, and another `<alphanum>` variable.

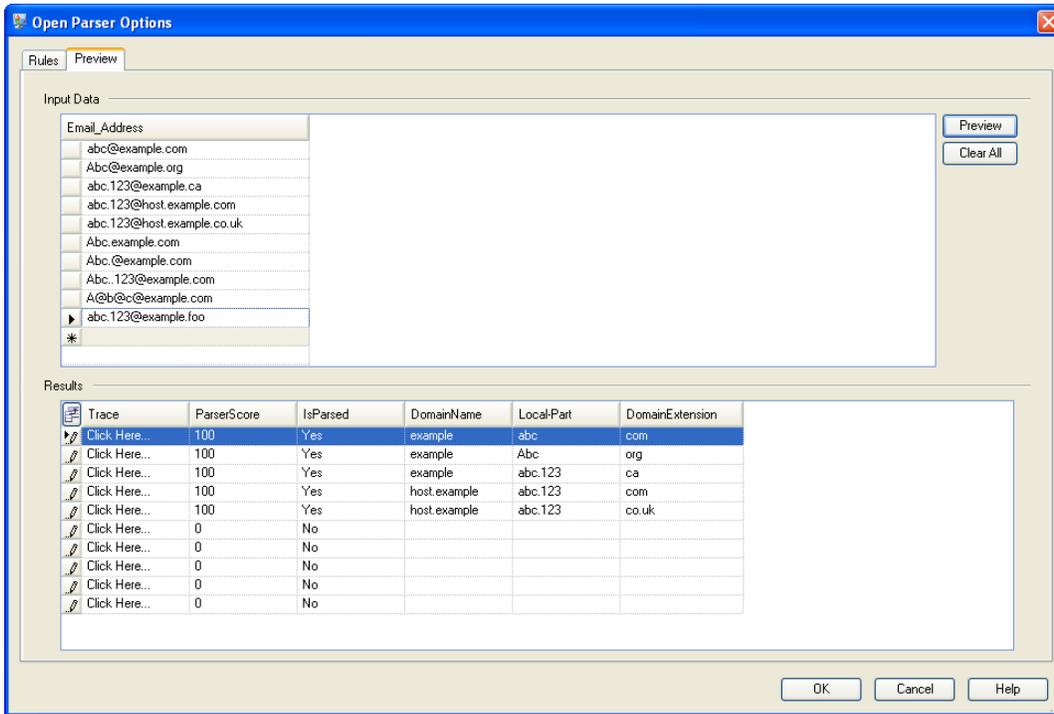
The `<alphanum>` variable definition is a regular expression that means any string of characters from A to Z, a to z, and 0-9. The `<alphanum>` variable is used throughout this parsing grammar and is defined once on the last line of the parsing grammar.

The parsing grammar uses a combination of regular expressions and literal characters to build a pattern for e-mail addresses. Any characters in double quotes in this parsing grammar are literal characters, the name of a table used for lookup, or a regular expression. The parsing grammar uses these special characters:

- The "+" character means that a regular expression can occur one or more times.
- The "?" character means that a regular expression can occur zero or one time.
- The "|" character means that the variable has an OR condition.
- The ";" character means end of a rule.

Use the **Commands** tab to explore the meaning of the other special symbols you can use in parsing grammars by hovering the mouse over the description.

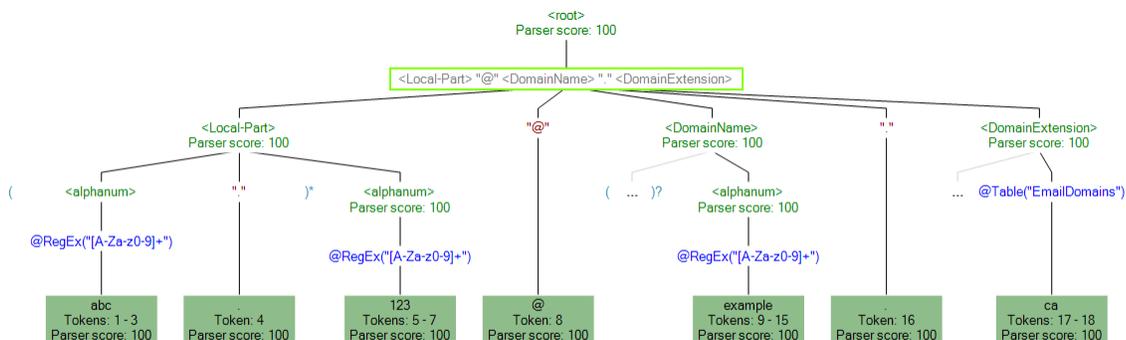
To test the parsing grammar, click the Preview tab. Type the e-mail addresses shown below in the **Email Address** field and then click **Preview**.



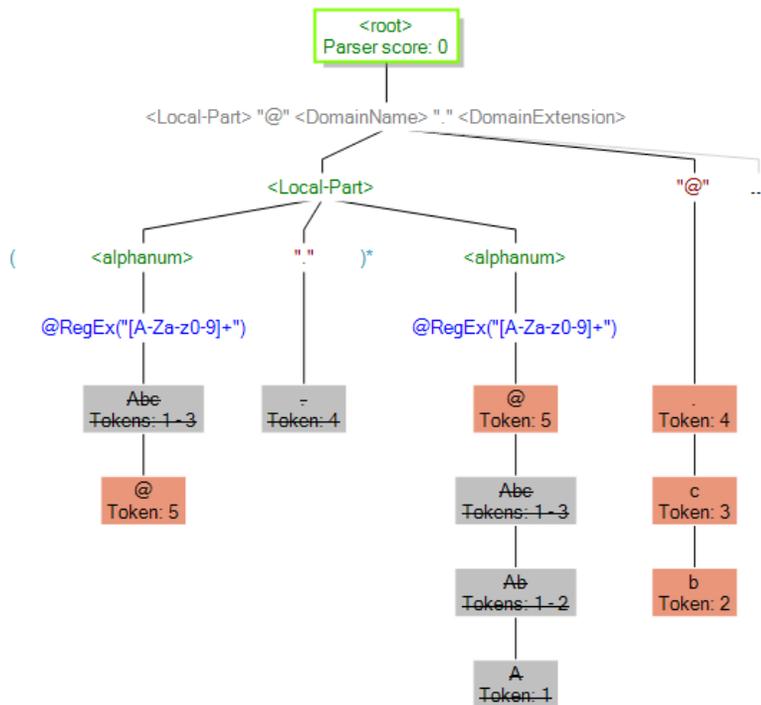
You can also type other e-mail addresses to see how the input data is parsed.

You can also use the Trace feature to see a graphical representation of either the final parsing results or to step through the parsing events. Click the link in the **Trace** column to see the Trace Details for the data row.

Trace Details shows a matching result. Compare the tokens matched for each expression in the parsing grammar.



You can also use Trace to view non-matching results. The following graphic shows a non-matching result. Compare the tokens matched for each expression in the parsing grammar. The reason that this input data (Abc.example.com) did not match is because it did not contain all of the required tokens to match—there is no @ character separating the Local-Part token and the Domain tokens.



Write to File

The template contains one Write to File stage. In addition to the input field, the output file contains the **Local-Part**, **DomainName**, **DomainExtension**, **IsParsed**, and **ParserScore** fields.

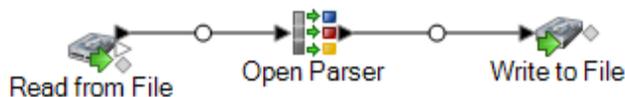
Parsing U.S. Phone Numbers

This template demonstrates how to parse U.S. phone numbers into component parts. The parsing rule separates each token in the **PhoneNumber** field and copies each token to four fields: **CountryCode**, **AreaCode**, **Exchange**, and **Number**.

Business Scenario

You work for a wireless provider and have been assigned a project to analyze incoming phone number data for a growing region of your business.

The following dataflow provides a solution to the business scenario:



This dataflow template is available in Enterprise Designer. Go to **File > New > Dataflow > From template** and select **ParseUSPhoneNumbers**. This dataflow requires Data Normalization.

In this dataflow, data is read from a file and processed through the Open Parser stage. For each data row in the input file, this data flow will do the following:

Read from File

This stage identifies the file name, location, and layout of the file that contains the phone numbers you want to parse.

Open Parser

This stage defines whether to use a culture-specific domain grammar created in the Domain Editor or to define a domain-independent grammar. A culture-specific parsing grammar that you create in the Domain Editor is a validated parsing grammar that is associated with a culture and a domain. A domain-independent parsing grammar that you create in Open Parser is a validated parsing grammar that is not associated with a culture and domain.

In this template, the parsing grammar is defined as a domain-independent grammar.

The Open Parser stage contains a parsing grammar that defines the following commands and expressions:

- %Tokenize is set to None. When Tokenize is set to None, the parsing grammar rule must include any spaces or other token separators within its rule definition.
- %InputField is set to parse input data from the **PhoneNumber** field.
- %OutputFields is set to separate parsed data into four fields: **CountryCode**, **AreaCode**, **Exchange**, and **Number**.
- The <root> expression defines pattern of tokens being parsed and includes OR statements (|), such that a valid phone number is:
 - **CountryCode**, **AreaCode**, **Exchange**, and **Number** OR
 - **AreaCode**, **Exchange**, and **Number** OR
 - **Exchange** and **Number**

The parsing grammar uses a combination of regular expressions and literal characters to build a pattern for phone numbers. Any characters in double quotes in this parsing grammar are literal characters or a regular expression.

The plus character (+) used in this <root> command is defined as a literal character because it is encapsulated in quotes. You can use single or double quotes to indicate a literal character. If the plus character is used without quotes, it means that the expression it follows can occur one or more times.

The phone number domain rules are defined to match the following character patterns:

- Zero or one occurrence of a "+" character.
- The **CountryCode** rule, which is a single digit between 0-9.

- Zero or one occurrence of an open parentheses or a hyphen or a space character. Two of these characters occurring in sequence results in a non-match, or in other words, an invalid phone number.
- The **AreaCode** rule, which is a sequence of exactly three digits between 0-9.
- Zero or one occurrence of an open parentheses or a hyphen or a space character. Two of these characters occurring in sequence results in a non-match, or in other words, an invalid phone number.
- The **Exchange** rule, which is a sequence of exactly three digits between 0-9.
- Zero or one occurrence of an open parentheses or a hyphen or a space character. Two of these characters occurring in sequence results in a non-match, or in other words, an invalid phone number.
- The **Number** rule, which is a sequence of exactly four digits between 0-9.

The rule variables that define the domain must use the same names as the output fields defined in the required `OutputFields` command.

Regular Expressions and Expression Quantifiers

The parsing grammar uses a combination of regular expressions and expression quantifiers to build a pattern for U.S. phone numbers. The parsing grammar uses these special characters:

- The "?" character means that a regular expression can occur zero or one time.
- The (|) character indicates an OR condition.
- The ";" character means end of a rule.

Use the **Commands** tab to explore the meaning of the other special symbols you can use in parsing grammars by hovering the mouse over the description.

Using the Preview Tab

To test the parsing grammar, click the Preview tab. Type the phone numbers shown below in the **PhoneNumber** field and then click **Preview**.

PhoneNumber	CountryCode	AreaCode	Exchange	Number
1(410)286-7334	1	410	286	7334
14042867534	1	404	286	7534
(410)286-7256		410	286	7256
301-868-9999		301	868	9999
1-222-458-7799	1	222	458	7799
+1(410)286-7334	1	410	286	7334
901 888 9990		901	888	9990
1 410 888 2345	1	410	888	2345
234-4567			234	4567
234 6789			234	6789

You can also type other valid and invalid phone numbers to see how the input data is parsed.

You can also use the Trace feature to see a graphical representation of either the final parsing results or to step through the parsing events. Click the link in the **Trace** column to see the Trace Details for the data row.

Write to File

The template contains one Write to File stage. In addition to the input field, the output file contains the **CountryCode**, **AreaCode**, **Exchange**, and **Number** fields.

3 - Standardization

In this section

Standardizing Terms.....	56
Standardizing Personal Names.....	57
Templates for Standardization.....	59



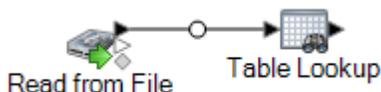
Standardizing Terms

Inconsistent use of terminology can be a data quality issue that causes difficulty in parsing, lookups, and more. You can create a dataflow that finds terms in your data that are inconsistently used and standardize them. For example, if your data includes the terms "Incorporated", "Inc.", and "Inc" in business names, you can create a dataflow to standardize on one form (for example, "Inc.").

Note: Before performing this procedure, your administrator must install the Data Normalization database containing standardized terms that you want to apply to your data. Instructions for installing databases can be found in the *Installation Guide*.

1. In Enterprise Designer, create a new dataflow.
2. Drag a source stage onto the canvas.
3. Double-click the source stage and configure it. See the *Dataflow Designer's Guide* for instructions on configuring source stages.
4. Drag a Table Lookup stage onto the canvas and connect it to the source stage.

For example, if you were using a Read from File source stage, your dataflow would look like this:



5. Double-click the Table Lookup stage on the canvas.
6. To specify the options for Table Lookup you create a rule. You can create multiple rules then specify the order in which you want to apply the rules. Click **Add** to create a rule.
7. In the **Action** field, leave the default option **Standardize** selected.
8. In the **On** field, leave **Complete field** selected if the whole field is the term you want to standardize. Or, choose **Individual terms within a field** to standardize individual words in the field.
9. In the **Source** field, select the field you want to standardize.
10. In the **Destination** field, select the field that you want to contain the standardized term. If you specify the same field as the source field, then the source field's value will be replaced with the standardized term.
11. In the **Table** field, select the table that contains the standardized terms.

Note: If you do not see the table you need, contact your system administrator. The Data Normalization database must be loaded.

12. In the **When table entry not found, set Destination's value to** field, select **Source's value**.
13. Click **OK**.

14. Define additional rules if you want to standardize values in more fields. When you are done defining rules, click **OK**.
15. Drag a sink stage onto the canvas and connect it to Table Lookup.

For example, if you were using Write to File, your dataflow would look like this:



16. Double-click the sink stage and configure it.

For information on configuring sink stages, see the *Dataflow Designer's Guide*.

You now have a dataflow that standardizes terms.

Standardizing Personal Names

This procedure shows how to create a dataflow that takes personal name data (for example "John P. Smith"), identifies common nicknames of the same name, and create a standard version of the name that can then be used to consolidate redundant records.

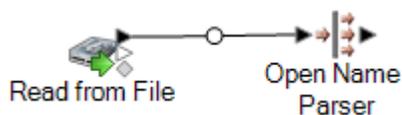
Note: Before beginning, make sure that your input data has a field named "Name" that contains the full name of the person.

1. If you have not already done so, load the following tables onto the Spectrum Technology Platform server:
 - Open Parser Base
 - Open Parser Enhanced Names

Use the Data Normalization database load utility to load these tables. For instructions on loading tables, see the *Installation Guide*.

2. In Enterprise Designer, create a new dataflow.
3. Drag a source stage onto the canvas.
4. Double-click the source stage and configure it. See the *Dataflow Designer's Guide* for instructions on configuring source stages.
5. Drag an Open Name Parser stage onto the canvas and connect it to the source stage.

For example, if you are using a Read from File stage, your dataflow would look like this:



6. Drag a Table Lookup stage onto the canvas and connect it to the Open Name Parser stage.
Your dataflow should now look like this:

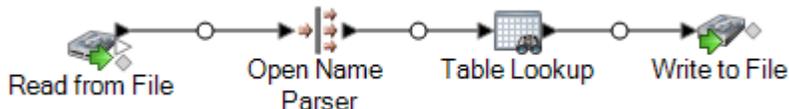


7. Double-click the Table Lookup stage on the canvas.
8. In the **Source** field, select **FirstName**.
9. In the **Destination** field, select **FirstName**.

By specifying the same field as both the source and destination, the field will be updated with the standardized version of the name.

10. In the **Table** field, select **NickNames.xml**.
11. Click **OK**.
12. Click **OK** again to close the **Table Lookup Options** window.
13. Drag a sink stage onto the canvas and connect it to the Table Lookup stage.

For example, if you were using a Write to File sink, your dataflow would now look like this:



14. Double-click the sink stage and configure it. See the *Dataflow Designer's Guide* for instructions on configuring source stages.

You now have a dataflow that takes personal names and standardizes the first name, replacing nicknames with the standard form of the name.

Templates for Standardization

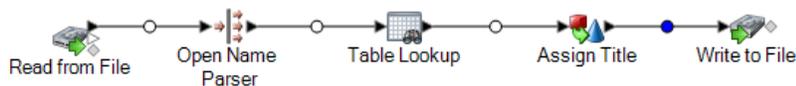
Formalizing Personal Names

This dataflow template demonstrates how to take personal name data (for example "John P. Smith"), identify common nicknames of the same name, and create a standard version of the name that can then be used to consolidate redundant records. It also show how you can add Title of Respect data based on Gender data.

Business Scenario

You work for a non-profit organization that wants to send out invitations for a gala event. Your input data include name data as full names and you want to parse the name data into First, Middle, and Last name fields and add a Title of Respect field to make your invitations more formal. You also want to replace any nicknames in your name data to use a more formal variant of the name.

The following dataflow provides a solution to the business scenario:



This dataflow template is available in Enterprise Designer. Go to **File > New > Dataflow > From template** and select **StandardizePersonalNames**. This dataflow requires the following products: Data Normalization and Universal Name.

For each data row in the input file, this data flow will do the following:

Read from File

This stage identifies the file name, location, and layout of the file that contains the names you want to parse. The file contains both male and female names.

Name Parser

In this template, the Name Parser stage is named Parse Personal Name. Parse Personal Name stage examines name fields and compares them to name data stored in the Spectrum Technology Platform name database files. Based on the comparison, it parses the name data into First, Middle, and Last name fields, assigns an entity type, and a gender to each name. It also uses pattern recognition in addition to the name data.

In this template the Parse Personal Name stage is configured as follows.

- Parse personal names is selected and Parse business names is cleared. When you select these options, first names are evaluated for gender, order, and punctuation and no evaluation of business names is performed.
- Gender Determination Source is set to default. For most cases, Default is the best setting for gender determination because it covers a wide variety of names. However, if you are processing names from a specific culture, select that culture. Selecting a specific culture helps ensure that the proper gender is assigned to the names. For example, if you leave Default selected, then the name Jean will be identified as a female name. However, if you select French, it will be identified as a male name.
- Order is set to natural. The name fields are ordered by Title, First Name, Middle Name, Last Name, and Suffix.
- Retain periods is cleared. Any punctuation in the name data is not retained.

Transformer

In this template, the Transformer stage is named Assign Titles. Assign Titles stage uses a custom script to search each row in the data stream output by the Parse Personal Name stage and assign a **TitleOfRespect** value based on the **GenderCode** value.

The custom script is:

```
if (row.get('TitleOfRespect') == '')
{
  if (row.get('GenderCode') == 'M')
    row.set('TitleOfRespect', 'Mr')
  if (row.get('GenderCode') == 'F')
    row.set('TitleOfRespect', 'Ms')
```

Every time the Assign Titles stage encounters **M** in the **GenderCode** field it sets the value for **TitleOfRespect** as **Mr**. Every time the Assign Titles stages encounters **F** in the **GenderCode** field it sets the value of **TitleOfRespect** as **Ms**.

Standardization

In this template, the Standardization stage is named Standardize Nicknames. Standardize Nickname stage looks up first names in the Nicknames.xml database and replaces any nicknames with the more regular form of the name. For example, the name Tommy is replaced with Thomas.

Write to File

The template contains one Write to File stage. In addition to the input fields, the output file contains the TitleOfRespect, FirstName, MiddleName, LastName, EntityType, GenderCode, and GenderDeterminationSource fields.

4 - Matching

In this section

Matching Terminology.....	62
Standard Fields	63
Techniques for Defining Match Keys.....	65
Match Rules.....	68
Matching Records from a Single Source.....	80
Matching Records from One Source to Another Source.....	86
Matching Records Between and Within Sources.....	92
Matching Records Against a Database.....	98
Matching Records Using Multiple Match Rules.....	100
Creating a Universal Matching Service.....	104
Using an Express Match Key.....	107
Analyzing Match Results.....	111
Dataflow Templates for Matching.....	126



Matching Terminology

Average Score	The average match score of all duplicates. The possible values are 0-100, with 0 indicating a poor match and 100 indicating an exact match.
Baseline	The selected match result that will be compared against another match result.
Candidate Group	Suspect and Candidate records grouped together by an ID assigned by CandidateFinder. The suspect (the first record in the group) is a record read from an Input source while its candidates are usually records found in a database using a SQL query.
Candidate Records	All non-suspect records in a match group or candidate group.
Drop	A decrease in duplicates.
Detail Match Record	A single record that corresponds to a record processed by a match stage. Each record provides information about whether the record was a Suspect, Unique, or a Duplicate as well as information about its Match Group or Candidate Group and output collection. Candidate records provide information on why the input record matched or did not match to its suspect.
Duplicate Collections	A duplicate collection consists of a Suspect and its Duplicate records grouped together by a CollectionNumber. Unique records always belong to CollectionNumber 0.
Duplicate Records	Number of records that match another record within a match group.
Express Matches	An express match is made when a suspect and candidate have an exact match on the contents of a designated field, usually an ExpressMatchKey provided by the Match Key Generator. If an Express Match is made no further processing is done to determine if the suspect and candidate are duplicates.
Input Records	Order of the records in the matching stage before the matching sort is performed.
Interflow Match	A matching stage that locates matches between similar data records between two input record streams. The first record stream is a source for suspect records and the second stream is a source for candidate records.
Intraflow Match	A matching stage that locates matches between similar data records within a single input stream.
Lift	An increase in duplicates.
Match Groups	(Group By) Records grouped together either by a match key or a sliding window.
Match Results	(or Resource Bundle) Logical grouping of files produced by a stage. This data is saved for each run of a stage and stored to disk. Subsequent runs will not overwrite or change the results from a previous run. In MAT, the bundles are used to provide information about the summary and details results, as well as settings information.

Match Results List	List of match results of a single type that MAT can analyze in the current analysis session.
Match Results Type	Indicates the contents of the match results. MAT uses the match results type to determine how to use the data.
Matcher Stage	A stage on the canvas that performs matching routines. The matcher stages are Interflow Match, Intraflow Match, and Transactional Match
Missed Match	A record that was previously a suspect or duplicate but is now unique.
New Match	A record that was previously unique but is now a suspect or duplicate.
Sliding Window	The sliding window matching method sequentially fills a predetermined buffer size called a window with the corresponding amount of data rows. As each row is added to the window it is compared to each item already contained in the window.
Suspect Records	A driver record that is matched against candidates within a match group or a candidate group.
Transactional Match	A matching stage that matches suspect records against candidate records that are returned from Candidate Finder or by an external application.
Unique Records	A suspect or candidate record that does not match any other records in a match group. If it is the only record in a match group, a suspect is automatically unique.

Standard Fields

This table lists some of the standard fields you come across in matching tasks, and gives description of these fields.

Fields	Descriptions
CandidateGroup	This field identifies a grouping of a suspect record and its candidates. Each suspect record is given a CandidateGroup number. The candidates for that suspect are given the same CandidateGroup number. For example, if John Smith is a suspect record and its candidate records are John Smith and Jon Smth, then all three records would have the same CandidateGroup value.
CollectionNumber	Identifies a collection of duplicate records. The possible values are 1 or values greater than 1. Records flagged with same value are duplicate records.

Fields	Descriptions
CollectionRecordType	<p>Identifies the template and best of breed records in a collection of duplicate records. The possible values are:</p> <ul style="list-style-type: none"> • Primary: The record is the selected template record in a collection. • Secondary: The record is not the selected template record in a collection. • BestOfBreed: The record is the newly created best of breed record in the collection.
ExpressMatchIdentified	<p>Indicates whether the match was obtained using the express match key. The possible values are <i>Yes</i> or <i>No</i>.</p>
HasDuplicates	<p>Identifies whether the record is a duplicate of another record. The options are:</p> <ul style="list-style-type: none"> • Y : The record is a suspect record and has duplicates. • N : The record is a suspect record and has no duplicates. • D : The record is a candidate record and is a duplicate of the suspect record. • U : The record is a candidate record but is not a duplicate of the suspect record.
MatchRecordType	<p>Identifies the type of match record in a collection. The possible values are:</p> <ul style="list-style-type: none"> • Suspect : A record that other records are compared to in order to determine if those are duplicates of each other. Each collection has one and only one suspect record. • Duplicate : A record that is a duplicate of the suspect record. • Unique : A record that has no duplicates.
MatchScore	<p>Identifies the overall score between two records. The possible values are 0-100, with 0 indicating a poor match and 100 indicating an exact match. Scores are generated basis the match rule is configured.</p>
TotalMatchCount	<p>This field indicates the total number of matches that were made during processing.</p>
TransactionRecordType	<p>It can be one of these:</p> <ul style="list-style-type: none"> • Suspect: A suspect record is used as input to a query. • Candidate: A candidate record is a result returned from a query.

Techniques for Defining Match Keys

Effective and efficient matching requires the right balance between accuracy and performance. The most accurate approach to matching would be to analyze each record against all other records, but this is not practical because the number of records that would need to be processed would result in unacceptably slow performance. A better approach is to limit the number of records involved in the matching process to those that are most likely to match. You can do this by using match keys. A match key is a value created for each record using an algorithm that you define. The algorithm takes values from the record and uses it to produce a match key value, which is stored as a new field in the record.

For example, if the incoming record is:

First Name - Fred
 Last Name - Mertz
 Postal Code - 21114-1687
 Gender Code - M

And you define a match key rule that generates a match key by combining data from the record like this:

Input Field	Start Position	Length
Postal Code	1	5
Postal Code	7	4
Last Name	1	5
First Name	1	5
Gender Code	1	1

Then the key would be:

211141687MertzFredM

Any records that have the same match key are placed into a match group. The matching process then compares records in the group to each other to identify matches.

To create a match key, use a Match Key Generator stage if you are matching records using Interflow Match or Intraflow Match. If you are matching records using Transactional Match, use the Candidate Finder stage to create match groups.

Note: The guidelines that follow can be applied to both Match Key Generator keys and Candidate Finder queries. In Candidate Finder, these guidelines apply to how you define the SELECT statement.

Match Group Size and Performance

The match key determines the size of the match group, and thus the performance of your dataflow. As the size of the match group doubles, execution time doubles. For example, if you define a match key that produces a group of 20 potentially-matching records, it will take twice as long to process as if you modify the match key so that the match group contains only 10 potentially-matching records. The disadvantage to "tightening" the match key rule to produce a smaller match group is that you run the risk of excluding records that do match. "Loosening" the match key rules reduces the chance of a matching record being excluded from the group, but increases group size. To find the right balance for your data it is important that you test with a variety of match key rules using a data that is representative of the data you intend to process in production.

Density

When designing a match key it is important to consider the density of the data. Density refers to the degree to which the data can be distributed across match groups. Since performance is determined by the number of comparisons the system has to perform, match keys that produce a small number of large match groups will result in slower performance than match keys that produce a large number of small match groups.

To illustrate this concept, consider a situation where you have a set of one million name and address records that you want to match. You might define a match key as the first three bytes of the postal code and the first letter of the last name. If the records are from all over the U.S., the match key would produce a good number of match groups and is likely to have acceptable performance. But if all the records are from New York, the postal codes would all begin with "100" and you would end up with, at most, only 26 match groups. This would produce large match groups containing, on average, approximately 38,000 records.

You can calculate the maximum number of comparisons performed for each match group by using the following formula:

$$N * (N-1) / 2$$

Where N is the number of records in the match group.

So if you have 26 match groups containing 38,000 records each, the maximum number of comparisons performed would be approximately 18.7 billion. Here is how this number is calculated:

First, determine the maximum number of comparisons per match group:

$$38,000 * (38,000-1) / 2 = 721,981,000$$

Then, multiply this amount by the number of match groups:

$$721,981,000 * 26 = 18,771,506,000$$

If there were instead 100 unique values for the first 3 bytes of the postal code you would have 2,600 match groups containing an average of 380 records. In this case the maximum number of comparisons would be 187 million, which is 100 times fewer. So if the records are only from New York, you might consider using the first four or even five bytes of the postal code for the match key in order to produce more match groups and reduce the number of comparisons. You may miss a few matches but the tradeoff would be greatly reduced execution time.

In reality, a match key like the one used in this example will not result in match groups of equal size because of variations in the data. For example, there will be many more people whose last name starts with "S" than with "X". Because of this, you should focus your efforts on reducing the size of the largest match groups. A match group of 100,000 records is 10 times larger than a match group of 10,000 but it will require 100 times more comparisons and will take 100 times as long. For example, say you are using five bytes of postal code and six bytes of the AddressLine1 field for your match key. On the surface that seems like a fairly fine match key. The problem is with PO Box addresses. While most of the match groups may be of an acceptable size, there would be a few very large match groups with keys like 10002PO BOX that contain a very large number of records. To break up the large match groups you could modify your match key to include the first couple of digits of the PO box number.

Aligning the Match Key with the Match Rule

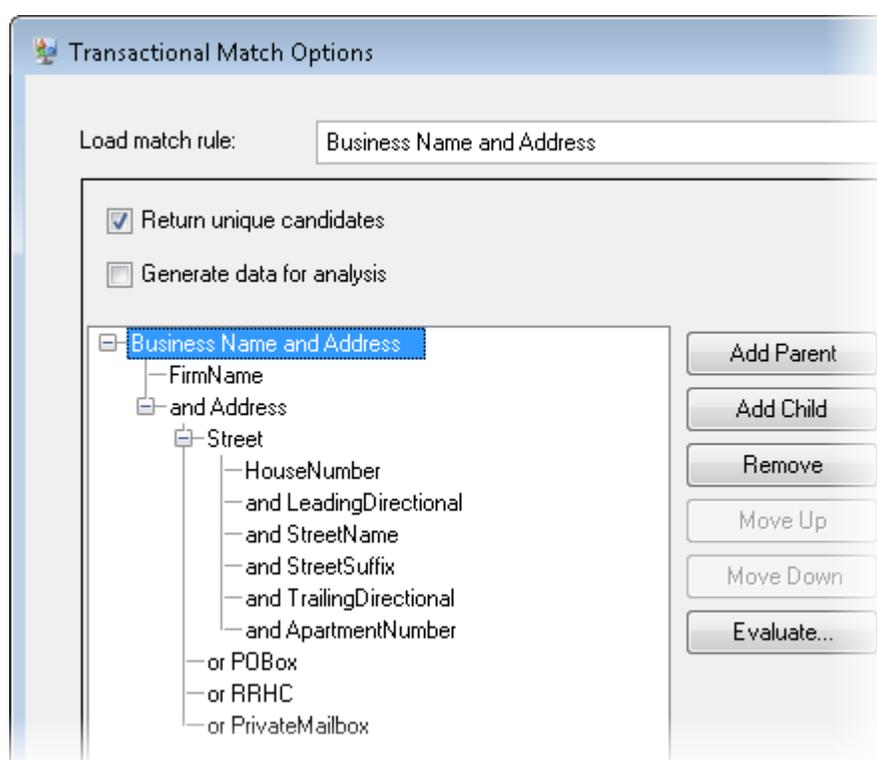
To achieve the most accurate results, you should design the match key to work well with the match rule that you will use it with. This requires you to consider how the match rule is defined.

- The match key should include any fields that the match rule requires to be an exact match.
- The match key should use the same kind of algorithm as is used in the match rule. For example, if you are designing a match key for use with a match rule that uses a phonetic algorithm, then the match key should also use a phonetic algorithm.
- The match key should be built using data from all the fields that are used in the match rule.
- Consider how the match key will be affected if there is data missing from one or more of the fields used for the match key. For example, say you use middle initial as part of the match key and you have a record for John A. Smith and another for John Smith. You have configured the match rule to ignore blank values in the middle initial field, so these two records would match according to your match rule. However, since the match key uses the middle initial, the two records would end up in different match groups and would not be compared to each other, thus defeating the intent of your match rule.

Match Rules

Each of the matching stages (Interflow Match, Intraflow Match, and Transactional Match) require you to configure a match rule. A match rule defines the criteria that are used to determine if one record matches another. It specifies the fields to compare, how to compare the fields, and a hierarchy of comparisons for complex matching rules.

Creating a hierarchical set of comparisons allows you to form nested Boolean match rules. For example, consider the following match rule:



In this example, the match rule is attempting to match records based on a business name and address. The first element of the match rule is the FirmName field. This element means that the value in the FirmName field must match in order for records to match. The second element evaluates the address. Note that it is prefaced with the logical operator "and" which means that both the FirmName and Address must match in order for records to match. The Address portion of the match rule consists of child rules that evaluate four types of addresses: street addresses, PO Box addresses, Rural Route/Highway Contract (RRHC) addresses, and private mailbox addresses. The Street child looks at the dataflow fields HouseNumber, LeadingDirectional, StreetName, StreetSuffix, TrailingDirectional, and ApartmentNumber. If all these match, then the parent rule "Street" and its parent rule "Address" all evaluate to "true". If the Street rule does not evaluate to true, the POBox

field is evaluated, then RRHC, then PrivateMailbox. If any of these three match then the parent Address element will match.

Building a Match Rule

Match rules are used in Interflow Match, Intraflow Match, and Transactional Match to define the criteria that determine if one record matches another. Match rules specify the fields to compare, how to compare the fields, and a hierarchy of comparisons for complex matching rules.

You can build match rules in Interflow Match, Intraflow Match, and Transactional Match. You can also build match rules in the Enterprise Designer Match Rule Management tool. Building a rule in the Match Rule Management tool makes the rule available to use in any dataflow, and also makes it available to other users. Building a match rule in one of the matcher stages makes the rule available only for that stage, unless you save the rule by clicking the **Save** button, which makes it available to other stages and users.

1. Open Enterprise Designer.
2. Do one of the following:
 - If you want to define a match rule in Interflow Match, Intraflow Match, or Transactional Match, double-click the match stage for which you want to define a match rule. In the **Load match rule** field, choose a predefined match rule as a starting point. If you want to start with a blank match rule, click **New**.
 - If you want to define a match rule in the Match Rule Management tool, select **Tools > Match Rule Management**. If you want to use an existing rule as a starting point for your rule, check the **Copy from** box and select the rule to use as a starting point.
3. Specify the dataflow fields you want to use in the match rule as well as the match rule hierarchy.
 - a) Click **Add Parent**.
 - b) Type in a name for the parent. The name must be unique and it cannot be a field. The first parent in the hierarchy is used as the match rule name in the **Load match rule** field. All custom match rules that you create and predefined rules that you modify are saved with the word "Custom" prepended to the name.
 - c) Click **Add Child**. A drop-down menu appears in the rule hierarchy. Select a field to add to the parent.

Note: All children under a parent must use the same logical operator. If you want to use different logical operators between fields you must first create intermediate parents.
 - d) Repeat to complete your matching hierarchy.
4. Define parent options. Parent options are displayed to the right of the rule hierarchy when a parent node is selected.

- a) Click **Match when not true** to change the logical operator for the parent from AND to AND NOT. If you select this option, records will only match if they do not match the logic defined in this parent.

Note: Checking the **Match when not true** option has the effect of negating the **Matching Method** options. For more information, see [Negative Match Conditions](#) on page 77.

- b) In the **Matching Method** field, specify how to determine if a parent is a match or a non-match. One of the following:

All true A parent is considered a match if all children are determined to match. This method creates an "AND" connector between children.

Any true A parent is considered a match if at least one child is determined to match. This method creates an "OR" connector between children.

Based on threshold A parent is considered a match if the score of the parent is greater than or equal to the parent's threshold. When you select this option, the **Threshold** slider appears. Use this slider to specify a threshold. The scoring method determines which logical connector to use. Thresholds at the parent cannot be higher than the threshold of the children.

Note: The threshold set here can be overridden at runtime in the Dataflow Options dialog box. Go to **Edit > Dataflow Options** and click **Add**. Expand the stage, click **Top level threshold**, and enter the threshold in the **Default value** field.

- c) In the **Missing Data** field, specify how to score blank data in a field. One of the following:

Ignore blanks Ignores the field if it contains blank data.

Count as 0 Scores the field as 0 if it contains blank data.

Count as 100 Scores the field as 100 if it contains blank data.

Compare Blanks Scores the suspect and candidate fields as 100 if they both contain blank data; otherwise, scores the suspect and candidate fields as 0.

- d) In the **Scoring method** field, select the method used for determining the matching score. One of the following:

Weighted Average Uses the weight of each child to determine the average match score.

Average Uses the average score of each child to determine the score of a parent.

Maximum Uses the highest child score to determine the score of a parent.

Minimum Uses the lowest child score to determine the score of a parent.

Vector Summation Uses the vector summation of each child score to determine the score of the parent. The formula for calculation is:

$\sqrt{a^2 + b^2 + c^2} / \sqrt{n}$, where: a, b, and c are the scores of three children and n is the number of children.

The following table shows the logical relationship between matching methods and scoring methods and how each combination changes the logic used during match processing.

Table 1: Matching Method-to-Scoring Method Matrix

Scoring Method	Matching Method			Comments
	Any True	All True	Based on Threshold	
Weighted Average	n/a	AND	AND	Only available when All True or Based on Threshold are selected as the Matching Method.
Average	n/a	AND	AND	
Vector Summation	n/a	AND	AND	
Maximum	OR	n/a	OR	Only available when Any True or Based on Threshold are selected as the Matching Method.
Minimum	OR	n/a	OR	

5. Define child options. Child options are displayed to the right of the rule hierarchy when a child is selected.
 - a) Check the option **Candidate field** to map the child record field selected to a field in the input file.
 - b) Check the option **Cross match against** and select one or more items from the dropdown list to match different fields to one another between two records. If you are using the Match Rule Management tool to create or edit a match rule, there will be no dropdown and you will instead need to enter each field name, separated by commas.
 - c) Click **Match when not true** to change the logical operator from AND to NOT. If you select this option, the match rule will only evaluate to true if the records do not match the logic defined in this child.

For example, if you want to identify individuals who are associated with multiple accounts, you could create a match rule that matches on name but where the account number does not match. You would use the **Match when not true** option for the child that matches the account number.

- d) In the **Missing Data** field, specify how to score blank data in a field. One of the following:
- | | |
|-----------------------|--|
| Ignore blanks | Ignores the field if it contains blank data. |
| Count as 0 | Scores the field as 0 if it contains blank data. |
| Count as 100 | Scores the field as 100 if it contains blank data. |
| Compare Blanks | Scores the suspect and candidate fields as 100 if they both contain blank data; otherwise, scores the suspect and candidate fields as 0. |
- e) In the **Threshold field**, specify the threshold that must be met at the individual field level in order for that field to be determined a match.
- f) In the **Scoring method** field, select the method used for determining the matching score. One of the following:
- | | |
|-------------------------|--|
| Weighted Average | Uses the weight of each algorithm to determine the average match score. |
| Average | Uses the average score of each algorithm to determine the match score. |
| Maximum | Uses the highest algorithm score to determine the match score. |
| Minimum | Uses the lowest algorithm score to determine the match score. |
| Vector Summation | Uses vector summation of the score of each algorithm to determine the match score. This scoring method is useful if you want a higher match score in one or more algorithms to get proportionately represented in the final match score. The formula used for calculating the final score is:

$\sqrt{a^2 + b^2 + c^2} / \sqrt{n}$ where: a, b, and c are the scores of three different algorithms and n is the number of algorithms used. |
- g) Choose one or more algorithms to use to determine if the values in the field match. One of the following.
- | | |
|--------------------------------|---|
| Acronym | Determines whether a business name matches its acronym by looking for acronym data; otherwise it creates an acronym using the first character of every word. Example: Internal Revenue Service and its acronym IRS would be considered a match and return a match score of 100. |
| Character Frequency | Determines the frequency of occurrence of each character in a string and compares the overall frequencies between two strings. |
| Daitch-Mokotoff Soundex | Phonetic algorithm that allows greater accuracy in matching of Slavic and Yiddish surnames with similar pronunciation but differences in spelling. Coded names are six digits long, and multiple possible encodings can be returned for a single name. This option was developed to respond to limitations of Soundex in the processing of Germanic or Slavic surnames. |

- Date** Compare date fields regardless of the date format in the input records. Click Edit in the Options column to specify the following:
- **Require Month:** prevents a date that consists only of a year from matching
 - **Require Day:** prevents a date that consists only of a month and year from matching
 - **Match Transposed MM/DD:** where month and day are provided in numeric format, compares suspect month to candidate day and suspect day to candidate month as well as the standard comparison of suspect month to candidate month and suspect day to candidate day
 - **Prefer DD/MM/YYYY format over MM/DD/YYYY:** contributes to date parsing in cases where both month and day are provided in numeric format and their identification can not be determined by context. For example, given the numbers 5 and 13, the parser will automatically assign 5 to the month and 13 to the day because there are only 12 months in a year. However, given the numbers 5 and 12 (or any two numbers 12 and under), the parser will assume whichever number is first to be the month. Checking this option will ensure that the parser reads the first number as the day rather than the month.
 - **Range Options—Overall:** allows you to set the maximum number of days between matching dates. For example, if you enter an overall range of 35 days and your candidate date is December 31st, 2000, a suspect date of February 5, 2001 would be a match, but a suspect date of February 6 would not. If you enter an overall range of 1 day and your candidate date is January 2000, a suspect date of 1999 would be a match (comparing December 31, 1999) but a suspect date of January 2001 would not.
 - **Range Options—Year:** allows you to set the number of years between matching dates, independent of month and day. For example, if you enter a year range of 3 and your candidate date is January 31, 2000, a suspect date of January 31, 2003, would be a match but a suspect date of February 2003 would not. Similarly, if your candidate date is 2000, a suspect date of March 2003 would be a match because months are not in conflict and it's within the three-year range.
 - **Range Options—Month:** allows you to set the number of months between matching dates, independent of year and day. For example, if you enter a month range of 4 and your candidate date is January 1, 2000, a suspect date of May 2000 is a match because there is no day conflict and it's within the four-month range, but a suspect date of May 2, 2000, is not, because the days conflict.
 - **Range Options—Day:** allows you to set the number of days between matching dates, independent of year and month. For example, if you enter a day range of 5 and your candidate date is January 1, 2000, a suspect date of January 2000 is a match because there is no day conflict

but a suspect date of December 27, 1999, is not, because the months conflict.

Double Metaphone	Determines the similarity between two strings based on a phonetic representation of their characters. Double Metaphone is an improved version of the Metaphone algorithm, and attempts to account for the many irregularities found in different languages.
Edit Distance	Determines the similarity between two strings based on the number of deletions, insertions, or substitutions required to transform one string into another.
Euclidean Distance	Provides a similarity measure between two strings using the vector space of combined terms as the dimensions. It also determines the greatest common divisor of two integers. It takes a pair of positive integers and forms a new pair that consists of the smaller number and the difference between the larger and smaller numbers. The process repeats until the numbers are equal. That number then is the greatest common divisor of the original pair. For example, 21 is the greatest common divisor of 252 and 105: ($252 = 12 \times 21$; $105 = 5 \times 21$); since $252 - 105 = (12 - 5) \times 21 = 147$, the GCD of 147 and 105 is also 21.
Exact Match	Determines if two strings are the same.
Initials	Used to match initials for parsed personal names.
Jaro-Winkler Distance	Determines the similarity between two strings based on the number of character replacements it takes to transform one string into another. This option was developed for short strings, such as personal names.
Keyboard Distance	Determines the similarity between two strings based on the number of deletions, insertions, or substitutions required to transform one string to the other, weighted by the position of the keys on the keyboard. Click Edit in the Options column to specify the type of keyboard you are using: QWERTY (U.S.), QWERTZ (Austria and Germany), or AZERTY (France).
Koeln	Indexes names by sound as they are pronounced in German. Allows names with the same pronunciation to be encoded to the same representation so that they can be matched, despite minor differences in spelling. The result is always a sequence of numbers; special characters and white spaces are ignored. This option was developed to respond to limitations of Soundex.
Kullback-Liebler Distance	Determines the similarity between two strings based on the differences between the distribution of words in the two strings.
Metaphone	Determines the similarity between two English-language strings based on a phonetic representation of their characters. This option was developed to respond to limitations of Soundex.

Metaphone (Spanish)	Determines the similarity between two strings based on a phonetic representation of their characters. This option was developed to respond to limitations of Soundex.
Metaphone 3	Improves upon the Metaphone and Double Metaphone algorithms with more exact consonant and internal vowel settings that allow you to produce words or names more or less closely matched to search terms on a phonetic basis. Metaphone 3 increases the accuracy of phonetic encoding to 98%. This option was developed to respond to limitations of Soundex.
Name Variant	Determines whether two names are variants of each other. The algorithm returns a match score of 100 if two names are variations of each other, and a match score of 0 if two names are not variations of each other. For example, JOHN is a variation of JAKE and returns a match score of 100. JOHN is not a variant of HENRY and returns a match score of 0. Click Edit in the Options column to select Name Variant options. For more information, see Name Variant Finder on page 366.
NGram Distance	Calculates in text or speech the probability of the next term based on the previous n terms, which can include phonemes, syllables, letters, words, or base pairs and can consist of any combination of letters. This algorithm includes an option to enter the size of the NGram; the default is 2.
NGram Similarity	Determines similarity between two strings based on the length of the longest common subsequence of phonemes, syllables, letters, words or base pairs. The algorithm includes the following options: <ul style="list-style-type: none"> • Ngram size: Enter the size of the NGram. The default value is 2. • Drop Noise Characters: Select the check-box to replace punctuation with space. • Drop Spaces: Select the check-box to merge words.
Numeric String	Compares address lines by separating the numerical attributes of an address line from the characters. For example, in the string address 1234 Main Street Apt 567, the numerical attributes of the string (1234567) are parsed and handled differently from the remaining string value (Main Street Apt). The algorithm first matches numeric data in the string with the numeric algorithm. If the numeric data match is 100, the alphabetic data is matched using Edit distance and Character Frequency. The final match score is calculated as follows: $(\text{numericScore} + (\text{EditDistanceScore} + \text{CharacterFrequencyScore}) / 2) / 2$ For example, the match score of these two addresses is 95.5, calculated as follows:

123 Main St Apt 567
123 Maon St Apt 567

Numeric Score = 100
Edit Distance = 91
Character Frequency = 91

$91 + 91 = 182$
 $182/2 = 91$
 $100 + 91 = 191$
 $191/2 = 95.5$

Nysiis	Phonetic code algorithm that matches an approximate pronunciation to an exact spelling and indexes words that are pronounced similarly. Part of the New York State Identification and Intelligence System. Say, for example, that you are looking for someone's information in a database of people. You believe that the person's name sounds like "John Smith", but it is in fact spelled "Jon Smath". If you conducted a search looking for an exact match for "John Smith" no results would be returned. However, if you index the database using the NYSIIS algorithm and search using the NYSIIS algorithm again, the correct match will be returned because both "John Smith" and "Jon Smath" are indexed as "JANSNATH" by the algorithm. This option was developed to respond to limitations of Soundex; it handles some multi-character n-grams and maintains relative vowel positioning, whereas Soundex does not. Note: This algorithm does not process non-alpha characters; records containing them will fail during processing.
Phonix	Preprocesses name strings by applying more than 100 transformation rules to single characters or to sequences of several characters. 19 of those rules are applied only if the character(s) are at the beginning of the string, while 12 of the rules are applied only if they are at the middle of the string, and 28 of the rules are applied only if they are at the end of the string. The transformed name string is encoded into a code that is comprised by a starting letter followed by three digits (removing zeros and duplicate numbers). This option was developed to respond to limitations of Soundex; it is more complex and therefore slower than Soundex.
Sonnex	This algorithm determines the similarity between two French-language strings based on the phonetic representation of their characters. It returns a Sonnex coded key of the selected fields.
Soundex	Determines the similarity between two strings based on a phonetic representation of their characters.
SubString	Determines whether one string occurs within another.

Syllable Alignment	Combines phonetic information with edit distance-based calculations. Converts the strings to be compared into their corresponding sequences of syllables and calculates the number of edits required to convert one sequence of syllables to the other.
---------------------------	---

The following table describes the logical relationship between the number of algorithms you can use based on the parent scoring method selected.

Table 2: Matching Algorithm-to-Scoring Method Matrix

Scoring Method	Algorithms	
	Single	Multiple
Weighted Average	n/a	Yes
Average	n/a	Yes
Maximum	Yes	Yes
Minimum	n/a	Yes
Vector Summation	n/a	Yes

- If you are defining a rule in Interflow Match, Intraflow Match, or Transactional Match, and you want to share the rule with other stages and/or users, click the **Save** button at the top of the window.

Negative Match Conditions

Match conditions are statements that indicate which fields you want to match in order for two records to be considered a match. However, in some situations you may want to define a condition that says that two fields must *not* match in order for two records to be considered a match. This technique, known as *negation*, reverses the logic of a condition within a match rule.

For example, say you have customer support records for a call center and you want to identify customers who have contacted the call center but done so for multiple accounts. In other words, you want to identify individuals who are associated with multiple accounts. In order to identify customers who have multiple accounts, you would want to match records where the name matches but the account number does *not* match. In this case you would use negation on a match condition for the account number.

To use negation, check the box **Match when not true** when defining your match rule. This option is available to both parents (groups of conditions) and children (individual conditions) in the match

rule. The effect of this option is slightly different when used on a parent as opposed to a child. When used on a parent, the **Match when not true** option effectively reverses the matching method option as follows:

- The **All true** matching method effectively becomes "any false". The match rule can only match records if at least one of the children under the parent evaluates to false, thus making the parent evaluate to false. Since the **Match when not true** option is enabled, this evaluation to false will result in a match.
- The **Any true** matching method effectively becomes "none true". The match rule can only match records where none of the children evaluate to true because if any of the children evaluate to true, the parent will be true, but with the **Match when not true** option enabled, this evaluation to true will not result in a match. Only if none of the children are true, resulting in the parent evaluating to "not true", can the rule find a match.
- The **Based on threshold** matching method effectively changes from matching records that are equal to or greater than a specified threshold, to matching records that are less than the threshold. This is because records with a threshold value less than the one specified will evaluate to false, and since **Match when not true** is enabled, this will result in a match.

The **Match when not true** option is easier to understand when applied to child elements in a match rule. It simply indicates that two records are considered a match if the algorithm does *not* indicate a match.

Testing a Match Rule

After defining a match rule you may want to test it to see its results. To do this, you can use Match Rule Evaluation to examine the effects of a match rule on a small set of sample data.

1. Open the dataflow in Enterprise Designer.
2. Double-click the stage containing the match rule you want to test.

Match rules are used in Interflow Match, Intraflow Match, and Transactional Match.

3. In the match rule hierarchy, select the node you want to test and click **Evaluate**.
4. On the **Import** tab, enter the test data (a suspect and up to 10 candidates). There are two ways to enter test data.
 - To type in the test data manually, type a suspect record under **Suspect** and up to ten candidates under **Candidate**. After typing the records, you can click **Export** to save the records to a file which you can import later instead of reentering the data manually.
 - To import test data from a file, click **Import...** and select the file containing the sample records. Delimited files can be comma, pipe or tab delimited and should have a header record with header fields that match the field names shown under **Candidates**. A sample header record for Household input would be:

```
Name,AddressLine1,City,StateProvince
```

5. Evaluate the rule using one of these methods:

- Click **Current Rule**. This runs the rule defined on the **Match Rule** tab. Results are displayed for one suspect and candidate pair at a time. To cycle through the results, click the arrow buttons. Scores for fields and algorithms are displayed in a tree format similar to the match rule control. The results can optionally be exported to an XML file.

Note: If you make changes to the match rule and want to apply the changes to the stage's match rule, click **Save**.

- Click **All Algorithms**. This ignores the match rule and instead runs all algorithms against each field for suspect and candidate pairs. Results are displayed for one suspect and candidate pair at a time and can be cycled through using the arrow buttons.

To automatically update the results as you make changes to the match rule and/or input, select the **Auto update** check box. When using this feature with the **All Algorithms** option, only changes to the input will update the results.

The results shown under Scores are color coded as follows:

- Green—The rule resulted in a match.
- Red—The rule that did not result in a match.
- Gray—The rule was ignored
- Blue—The results for individual algorithms within the rule.

To export the evaluation results in XML format, click **Export**.

Sharing a Match Rule

You can create match rules that can be shared among modules, stages, dataflows, and users. By sharing a match rule, you can make it easier to develop dataflows by defining a match rule once and then referencing it wherever needed. This also helps to ensure that the match rules intended to perform the same functions are consistent across dataflows.

The match rules saved in the Match Rule Repository can be exported as .mr or .json files using the Administration utility. This is especially useful when trying to share these rules with other Spectrum installations or while consuming these in the Spectrum Big Data Quality SDK.

Note: By default the match rules are exported as .mr file.

- To share a match rule you built in an Interflow Match, Intraflow Match, or Transactional Match, click the **Save** button at the top of the stage's options window.
- If you build the rule in the Match Rules Management tool, the rule is automatically available to use in dataflows by all users. To view the Match Rules Management tool, in Enterprise Designer select **Tools > Match Rules Management**.

Viewing Shared Match Rules

In Enterprise Designer you can browse all the shared match rules available on your Spectrum Technology Platform system. These match rules can be used by Interflow Match, Intraflow Match, and Transactional Match stages in a dataflow to perform matching.

To browse the match rules in the Match Rule Repository, follow this procedure.

1. Open Enterprise Designer.
2. Select **Tools > Match Rules Management**.
3. Select the rule you want to view and click **View**.

Creating a Custom Match Rule as a JSON Object

Match rules can be configured and passed at runtime if they are exposed as dataflow options. This enables you to share match rules across machines and override existing match rules with JSON-formatted match rule strings. You can also set stage options when calling the job through a process flow or through the job executor command-line tool.

You can find schemas for MatchRule and MatchInfo in the following folder:

```
<Spectrum Location>\server\modules\jsonSchemas\matcher.
```

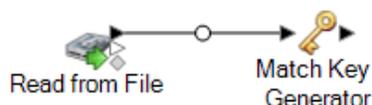
1. Save and expose the dataflow that contains the match rule.
2. Open the dataflow that uses the match rule.
3. Go to `Edit > Dataflow Options`.
4. In the **Map dataflow options to stages** table, click the matching stage that uses the match rule and check the **Custom Match Rule** box.
5. Optional: Change the name of the match rule in the **Option label** field from "Custom Match Rule" to the name you prefer.
6. Click **OK** twice.

Matching Records from a Single Source

This procedure describes how to use an Intraflow Match stage to identify groups of records within a single data source (such as a file or database table) that are related to each other based on the matching criteria you specify. The dataflow groups records into collections and writes the collections to an output file.

1. In Enterprise Designer, create a new dataflow.
2. Drag a source stage onto the canvas.
3. Double-click the source stage and configure it. See the *Dataflow Designer's Guide* for instructions on configuring source stages.
4. Drag a Match Key Generator stage onto the canvas and connect it to the source stage.

For example, if you are using a Read from File source stage, your dataflow would now look like this:



Match Key Generator creates a non-unique key for each record, which can then be used by matching stages to identify groups of potentially duplicate records. Match keys facilitate the matching process by allowing you to group records by match key and then only comparing records within these groups.

5. Double-click Match Key Generator.
6. Click **Add**.
7. Define the rule to use to generate a match key for each record.

Table 3: Match Key Generator Options

Option Name	Description and Valid Values
-------------	------------------------------

Algorithm	
-----------	--

Option Name	Description and Valid Values
-------------	------------------------------

Specifies one of these algorithms to use to generate the match key:

Consonant	Returns specified fields with consonants removed.
Double Metaphone	Returns a code based on a phonetic representation of their characters. Double Metaphone is an improved version of the Metaphone algorithm, and attempts to account for the many irregularities found in different languages.
Koeln	Indexes names by sound as they are pronounced in German. Allows names with the same pronunciation to be encoded to the same representation so that they can be matched, despite minor differences in spelling. The result is always a sequence of numbers; special characters and white spaces are ignored. This option was developed to respond to limitations of Soundex.
MD5	A message digest algorithm that produces a 128-bit hash value. This algorithm is commonly used to check data integrity.
Metaphone	Returns a Metaphone coded key of selected fields. Metaphone is an algorithm for coding words using their English pronunciation.
SpanishMetaphone	Returns a Metaphone coded key of selected fields for the Spanish language. This metaphone algorithm codes words using their Spanish pronunciation.
Metaphone 3	Improves upon the Metaphone and Double Metaphone algorithms with more exact consonant and internal vowel settings that allow you to produce words or names more or less closely matched to search terms on a phonetic basis. Metaphone 3 increases the accuracy of phonetic encoding to 98%. This option was developed to respond to limitations of Soundex.
Nysis	Phonetic code algorithm that matches an approximate pronunciation to an exact spelling and indexes words that are pronounced similarly. Part of the New York State Identification and Intelligence System. Say, for example, that you are looking for someone's information in a database of people. You believe that the person's name sounds like "John Smith", but it is in fact spelled "Jon Smyth". If you conducted a search looking for an exact match for "John Smith" no results would be returned. However, if you index the database using the NYSIIS algorithm and search using the NYSIIS algorithm again, the correct match will be returned because both "John Smith" and "Jon Smyth" are indexed as "JAN SNATH" by the algorithm.
Phonix	Preprocesses name strings by applying more than 100 transformation rules to single characters or to sequences of several characters. 19 of those rules are applied only if the characters are at the beginning of the string, while 12 of the rules are applied only if they are at the middle of the string, and 28 of the rules are applied only if they are at the end of

Option Name	Description and Valid Values
	<p>the string. The transformed name string is encoded into a code that is comprised by a starting letter followed by three digits (removing zeros and duplicate numbers). This option was developed to respond to limitations of Soundex; it is more complex and therefore slower than Soundex.</p>
	<p>Sonnex This algorithm determines the similarity between two French-language strings based on the phonetic representation of their characters. It returns a Sonnex coded key of the selected fields.</p>
	<p>Soundex Returns a Soundex code of selected fields. Soundex produces a fixed-length code based on the English pronunciation of a word.</p>
	<p>Substring Returns a specified portion of the selected field.</p>
Field name	<p>Specifies the field to which you want to apply the selected algorithm to generate the match key. For example, if you select a field called LastName and you choose the Soundex algorithm, the Soundex algorithm would be applied to the data in the LastName field to produce a match key.</p>
Start position	<p>Specifies the starting position within the specified field. Not all algorithms allow you to specify a start position.</p>
Length	<p>Specifies the length of characters to include from the starting position. Not all algorithms allow you to specify a length.</p>
Remove noise characters	<p>Removes all non-numeric and non-alpha characters such as hyphens, white space, and other special characters from an input field.</p>
Sort input	<p>Sorts all characters in an input field or all terms in an input field in alphabetical order.</p> <p>Characters Sorts the characters values from an input field prior to creating a unique ID.</p> <p>Terms Sorts each term value from an input field prior to creating a unique ID.</p>

- When you are done defining the rule click **OK**.

9. If you want to add additional match rules, click **Add** and add them, otherwise click **OK** when you are done.
10. Drag an Intraflow Match stage onto the canvas and connect it to the Match Key Generator stage.

For example, if you are using a Read from File source stage, your dataflow would now look like this:



11. Double-click Intraflow Match.
12. In the **Load match rule** field, select one of the predefined match rules which you can either use as-is or modify to suit your needs. If you want to create a new match rule without using one of the predefined match rules as a starting point, click **New**. You can only have one custom rule in a dataflow.

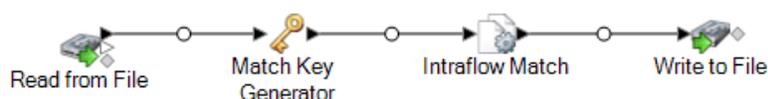
Note: The Dataflow Options feature in Enterprise Designer enables the match rule to be exposed for configuration at runtime.

13. In the **Group by** field, select **MatchKey**.

This will place records that have the same match key into a group. The match rule is applied to records within a group to see if there are duplicates. The match key for each record will be generated by the Generate Match Key stage you configured earlier in this procedure.

14. For information about modifying the other options, see [Building a Match Rule](#) on page 69.
15. Click **OK** to save your Intraflow Match configuration and return to the dataflow canvas.
16. Drag a sink stage onto the canvas and connect it to the Generate Match key stage.

For example, if you were using a Write to File sink stage your dataflow would look like this:



17. Double-click the sink stage and configure it.

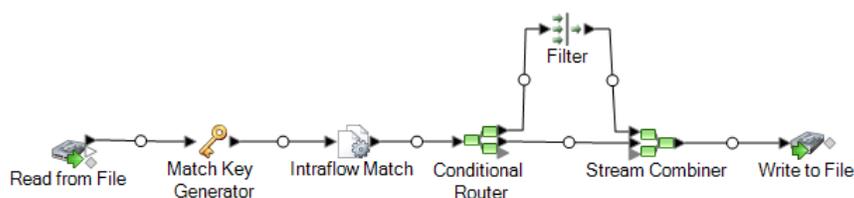
For information on configuring sink stages, see the *Dataflow Designer's Guide*.

You now have a dataflow that will match records from a single source.

Example of Matching Records in a Single Data Source

As a data steward for a credit card company, you want to analyze your customer database and find out which addresses occur multiple times and under what names so that you can minimize the number of duplicate credit card offers sent to the same household.

This example demonstrates how to identify members of the same household by comparing information within a single input file and creating an output file containing one record per household.



The Read from File stage reads in data that contains both unique records for each household and records that are potentially from the same household. The input file contains names and addresses.

The Match Key Generator creates a match key which is a non-unique key shared by like records that identify records as potential duplicates.

The Intraflow Match stage compares records that have the same match key and marks each record as either a unique record or as one of multiple records for the same household.

The Conditional Router sends records that are collections of records for each household to the Filter stage, which filters out all but one of the records from each household, and sends it on to the Stream Combiner stage. The Conditional Router stage also sends unique records directly to Stream Combiner.

Finally, the Write to File stage creates an output file that contains one record for each household.

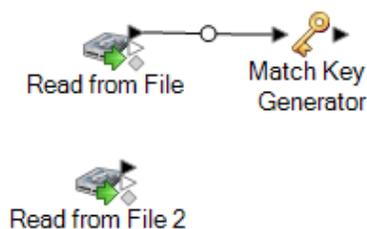
Matching Records from One Source to Another Source

This procedure describes how to use an Interflow Match stage to identify records in one source that match records in another source. The first source contains suspect records and the second source contains candidate records. The dataflow only matches records from one source to records in another source. It does not attempt to match records from within the same source. The dataflow groups records into collections of matching records and writes these collections to an output file.

1. In Enterprise Designer, create a new dataflow.
2. Drag two source stages onto the canvas. Configure one of them to point to the source of the suspect records and configure the other to point to the source of the candidate records.

See the *Dataflow Designer's Guide* for instructions on configuring source stages.

3. Drag a Match Key Generator stage onto the canvas and connect it to one of the source stages. For example, if you are using a Read from File source stage, your dataflow would now look like this:



Match Key Generator creates a non-unique key for each record, which can then be used by matching stages to identify groups of potentially duplicate records. Match keys facilitate the matching process by allowing you to group records by match key and then only comparing records within these groups.

Note: You will add a second Match Key Generator stage later. For now you only need one on the canvas.

4. Double-click the Match Key Generator stage.
5. Click **Add**.
6. Define the rule to use to generate a match key for each record.

Table 4: Match Key Generator Options

Option Name	Description and Valid Values
-------------	------------------------------

Algorithm	
-----------	--

Option Name	Description and Valid Values
-------------	------------------------------

Specifies one of these algorithms to use to generate the match key:

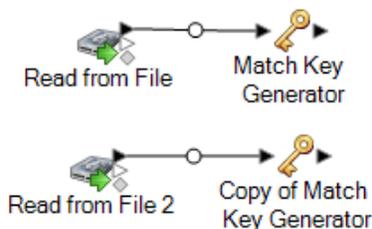
Consonant	Returns specified fields with consonants removed.
Double Metaphone	Returns a code based on a phonetic representation of their characters. Double Metaphone is an improved version of the Metaphone algorithm, and attempts to account for the many irregularities found in different languages.
Koeln	Indexes names by sound as they are pronounced in German. Allows names with the same pronunciation to be encoded to the same representation so that they can be matched, despite minor differences in spelling. The result is always a sequence of numbers; special characters and white spaces are ignored. This option was developed to respond to limitations of Soundex.
MD5	A message digest algorithm that produces a 128-bit hash value. This algorithm is commonly used to check data integrity.
Metaphone	Returns a Metaphone coded key of selected fields. Metaphone is an algorithm for coding words using their English pronunciation.
SpanishMetaphone	Returns a Metaphone coded key of selected fields for the Spanish language. This metaphone algorithm codes words using their Spanish pronunciation.
Metaphone 3	Improves upon the Metaphone and Double Metaphone algorithms with more exact consonant and internal vowel settings that allow you to produce words or names more or less closely matched to search terms on a phonetic basis. Metaphone 3 increases the accuracy of phonetic encoding to 98%. This option was developed to respond to limitations of Soundex.
Nysis	Phonetic code algorithm that matches an approximate pronunciation to an exact spelling and indexes words that are pronounced similarly. Part of the New York State Identification and Intelligence System. Say, for example, that you are looking for someone's information in a database of people. You believe that the person's name sounds like "John Smith", but it is in fact spelled "Jon Smyth". If you conducted a search looking for an exact match for "John Smith" no results would be returned. However, if you index the database using the NYSIIS algorithm and search using the NYSIIS algorithm again, the correct match will be returned because both "John Smith" and "Jon Smyth" are indexed as "JAN SNATH" by the algorithm.
Phonix	Preprocesses name strings by applying more than 100 transformation rules to single characters or to sequences of several characters. 19 of those rules are applied only if the characters are at the beginning of the string, while 12 of the rules are applied only if they are at the middle of the string, and 28 of the rules are applied only if they are at the end of

Option Name	Description and Valid Values
	<p>the string. The transformed name string is encoded into a code that is comprised by a starting letter followed by three digits (removing zeros and duplicate numbers). This option was developed to respond to limitations of Soundex; it is more complex and therefore slower than Soundex.</p>
	<p>Sonnex This algorithm determines the similarity between two French-language strings based on the phonetic representation of their characters. It returns a Sonnex coded key of the selected fields.</p>
	<p>Soundex Returns a Soundex code of selected fields. Soundex produces a fixed-length code based on the English pronunciation of a word.</p>
	<p>Substring Returns a specified portion of the selected field.</p>
Field name	<p>Specifies the field to which you want to apply the selected algorithm to generate the match key. For example, if you select a field called LastName and you choose the Soundex algorithm, the Soundex algorithm would be applied to the data in the LastName field to produce a match key.</p>
Start position	<p>Specifies the starting position within the specified field. Not all algorithms allow you to specify a start position.</p>
Length	<p>Specifies the length of characters to include from the starting position. Not all algorithms allow you to specify a length.</p>
Remove noise characters	<p>Removes all non-numeric and non-alpha characters such as hyphens, white space, and other special characters from an input field.</p>
Sort input	<p>Sorts all characters in an input field or all terms in an input field in alphabetical order.</p> <p>Characters Sorts the characters values from an input field prior to creating a unique ID.</p> <p>Terms Sorts each term value from an input field prior to creating a unique ID.</p>

7. When you are done defining the rule click **OK**.
8. Right-click the Match Key Generator stage on the canvas and select **Copy Stage**.

9. Right-click in an empty area of the canvas and select **Paste**.
10. Connect the copy of Match Key Generator to the other source stage.

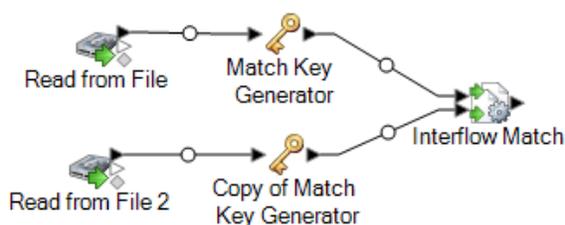
For example, if you are using Read from File input stages your dataflow would now look like this:



The dataflow now contains two Match Key Generator stages that produce match keys for each source using exactly the same rules. Having identically-configured Match Key Generator stages is essential to the proper functioning of this dataflow.

11. Drag an Interflow Match stage onto the canvas and connect each of the Match Key Generator stages to it.

For example, if you are using Read from File input stages your dataflow would now look like this:



12. Double-click the Interflow Match stage.
13. In the **Load match rule** field, select one of the predefined match rules which you can either use as-is or modify to suit your needs. If you want to create a new match rule without using one of the predefined match rules as a starting point, click **New**. You can only have one custom rule in a dataflow.

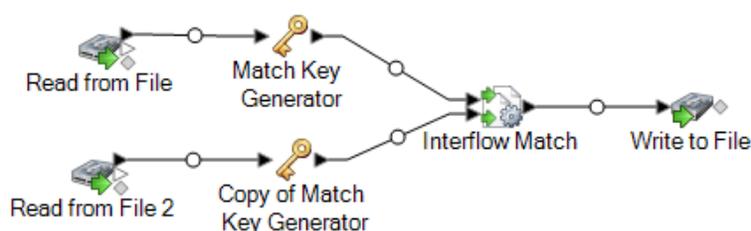
Note: The Dataflow Options feature in Enterprise Designer enables the match rule to be exposed for configuration at runtime.

14. In the **Group by** field, select **MatchKey**.

This will place records that have the same match key into a group. The match rule is applied to records within a group to see if there are duplicates. The match key for each record will be generated by the Generate Match Key stages you configured earlier in this procedure.

15. For information about modifying the other options, see [Building a Match Rule](#) on page 69.
16. Drag a sink stage onto the canvas and connect it to the Interflow Match stage.

For example, if you were using a Write to File sink stage your dataflow would look like this:



17. Double-click the sink stage and configure it.

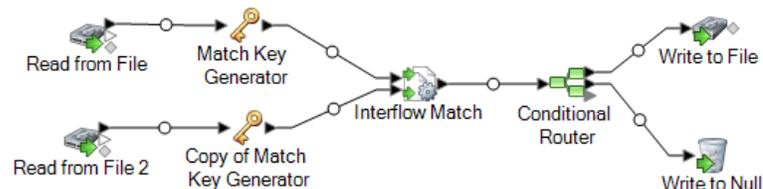
For information on configuring sink stages, see the *Dataflow Designer's Guide*.

You now have a dataflow that will match records from two data sources.

Example of Matching Records from Multiple Sources

As a direct mail company, you want to identify people who are on a do-not-mail list so that you do not send direct mail to them. You have a list of recipients in one file, and a list of people who do not wish to receive direct marketing mail in another file (a suppression file).

The following dataflow provides a solution to this business scenario:



The Read from File stage reads data from your mailing list, and the Read from File 2 stage reads data from the suppression list. The two Match Key Generator stages are identically configured so that they produce a match key which can be used by Interflow Match to form groups of potential matches. Interflow Match identifies records in the mailing list that are also in the suppression file and marks these records as duplicates. Conditional Router sends unique records, meaning those records that were not found in the suppression list, to Write to File to be written out to a file. The Conditional Router stage sends all other records to Write to Null where they are discarded.

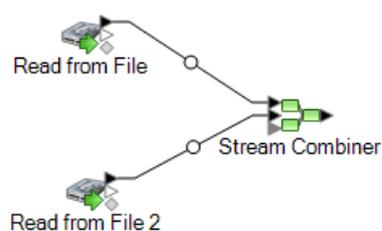
Matching Records Between and Within Sources

This procedure describes how to use an Intraflow Match stage to identify records in one file that match records in another file and in the same file. For example, you have two files (file A and file B) and you want to see if there are records in file A that match records in file B, but you also want to

see if there are records in file A that match other records in file A. You can accomplish this using a Stream Combiner and an Intraflow Match stage.

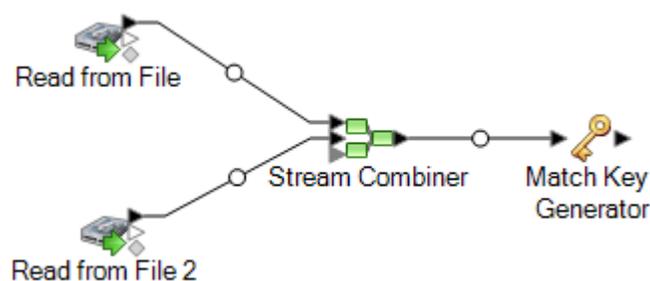
1. In Enterprise Designer, create a new dataflow.
2. Drag a source stage onto the canvas.
3. Double-click the source stage and configure it. See the *Dataflow Designer's Guide* for instructions on configuring source stages.
4. Drag a second source stage onto the canvas and configure it to read the second data source into the dataflow.
5. Drag a Stream Combiner stage onto the canvas and connect the two source stages to it.

For example, if your dataflow had two Read from File stages it would look like this after adding the Stream Combiner:



6. Drag a Match Key Generator stage onto the canvas and connect it to the Stream Combiner stage.

For example, your dataflow may now look like this:



Match Key Generator creates a non-unique key for each record, which can then be used by matching stages to identify groups of potentially duplicate records. Match keys facilitate the matching process by allowing you to group records by match key and then only comparing records within these groups.

7. Double-click Match Key Generator.
8. Click **Add**.
9. Define the rule to use to generate a match key for each record.

Table 5: Match Key Generator Options

Option Name	Description and Valid Values
-------------	------------------------------

Algorithm	
-----------	--

Option Name	Description and Valid Values
-------------	------------------------------

Specifies one of these algorithms to use to generate the match key:

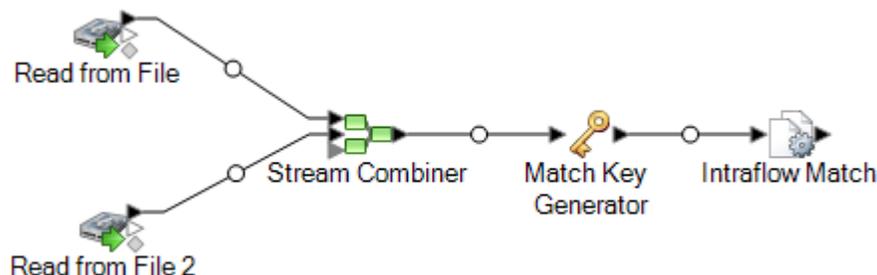
Consonant	Returns specified fields with consonants removed.
Double Metaphone	Returns a code based on a phonetic representation of their characters. Double Metaphone is an improved version of the Metaphone algorithm, and attempts to account for the many irregularities found in different languages.
Koeln	Indexes names by sound as they are pronounced in German. Allows names with the same pronunciation to be encoded to the same representation so that they can be matched, despite minor differences in spelling. The result is always a sequence of numbers; special characters and white spaces are ignored. This option was developed to respond to limitations of Soundex.
MD5	A message digest algorithm that produces a 128-bit hash value. This algorithm is commonly used to check data integrity.
Metaphone	Returns a Metaphone coded key of selected fields. Metaphone is an algorithm for coding words using their English pronunciation.
SpanishMetaphone	Returns a Metaphone coded key of selected fields for the Spanish language. This metaphone algorithm codes words using their Spanish pronunciation.
Metaphone 3	Improves upon the Metaphone and Double Metaphone algorithms with more exact consonant and internal vowel settings that allow you to produce words or names more or less closely matched to search terms on a phonetic basis. Metaphone 3 increases the accuracy of phonetic encoding to 98%. This option was developed to respond to limitations of Soundex.
Nysis	Phonetic code algorithm that matches an approximate pronunciation to an exact spelling and indexes words that are pronounced similarly. Part of the New York State Identification and Intelligence System. Say, for example, that you are looking for someone's information in a database of people. You believe that the person's name sounds like "John Smith", but it is in fact spelled "Jon Smyth". If you conducted a search looking for an exact match for "John Smith" no results would be returned. However, if you index the database using the NYSIIS algorithm and search using the NYSIIS algorithm again, the correct match will be returned because both "John Smith" and "Jon Smyth" are indexed as "JAN SNATH" by the algorithm.
Phonix	Preprocesses name strings by applying more than 100 transformation rules to single characters or to sequences of several characters. 19 of those rules are applied only if the characters are at the beginning of the string, while 12 of the rules are applied only if they are at the middle of the string, and 28 of the rules are applied only if they are at the end of

Option Name	Description and Valid Values
	<p>the string. The transformed name string is encoded into a code that is comprised by a starting letter followed by three digits (removing zeros and duplicate numbers). This option was developed to respond to limitations of Soundex; it is more complex and therefore slower than Soundex.</p>
	<p>Sonnex This algorithm determines the similarity between two French-language strings based on the phonetic representation of their characters. It returns a Sonnex coded key of the selected fields.</p>
	<p>Soundex Returns a Soundex code of selected fields. Soundex produces a fixed-length code based on the English pronunciation of a word.</p>
	<p>Substring Returns a specified portion of the selected field.</p>
Field name	<p>Specifies the field to which you want to apply the selected algorithm to generate the match key. For example, if you select a field called LastName and you choose the Soundex algorithm, the Soundex algorithm would be applied to the data in the LastName field to produce a match key.</p>
Start position	<p>Specifies the starting position within the specified field. Not all algorithms allow you to specify a start position.</p>
Length	<p>Specifies the length of characters to include from the starting position. Not all algorithms allow you to specify a length.</p>
Remove noise characters	<p>Removes all non-numeric and non-alpha characters such as hyphens, white space, and other special characters from an input field.</p>
Sort input	<p>Sorts all characters in an input field or all terms in an input field in alphabetical order.</p> <p>Characters Sorts the characters values from an input field prior to creating a unique ID.</p> <p>Terms Sorts each term value from an input field prior to creating a unique ID.</p>

10. When you are done defining the rule click **OK**.

11. If you want to add additional match rules, click **Add** and add them, otherwise click **OK** when you are done.
12. Drag an Intraflow Match stage onto the canvas and connect it to the Match Key Generator stage.

For example, your dataflow may now look like this:



13. Double-click Intraflow Match.
14. In the **Load match rule** field, select one of the predefined match rules which you can either use as-is or modify to suit your needs. If you want to create a new match rule without using one of the predefined match rules as a starting point, click **New**. You can only have one custom rule in a dataflow.

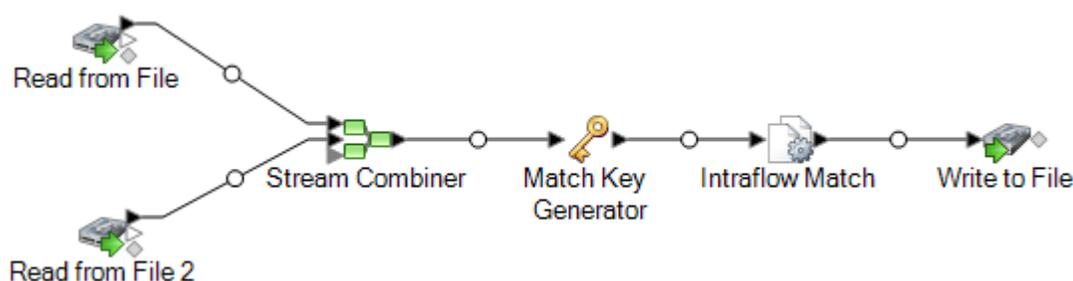
Note: The Dataflow Options feature in Enterprise Designer enables the match rule to be exposed for configuration at runtime.

15. In the **Group by** field, select **MatchKey**.

This will place records that have the same match key into a group. The match rule is applied to records within a group to see if there are duplicates. The match key for each record will be generated by the Generate Match Key stage you configured earlier in this procedure.

16. For information about modifying the other options, see [Building a Match Rule](#) on page 69.
17. Click **OK** to save your Intraflow Match configuration and return to the dataflow canvas.
18. Drag a sink stage onto the canvas and connect it to the Generate Match key stage.

For example, if you were using a Write to File sink stage your dataflow would look like this:



19. Double-click the sink stage and configure it.

For information on configuring sink stages, see the *Dataflow Designer's Guide*.

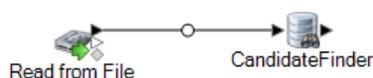
Matching Records Against a Database

This procedure describes how to match records where the suspect records come from a source such as a file or database, and the candidate records are in a database with other unrelated records. For each input record, the dataflow queries the database for candidates for that record, then uses a Transactional Match stage to match records. Finally, the dataflow writes the collections of matching records to an output file.

Note: Transactional Match only matches suspect records to candidates. It does not attempt to match suspect records to other suspect records as is done in Intraflow Match.

1. In Enterprise Designer, create a new dataflow.
2. Drag a source stage onto the canvas.
3. Double-click the source stage and configure it. See the *Dataflow Designer's Guide* for instructions on configuring source stages.
4. Drag a Candidate Finder stage to the canvas and connect the source stage to it.

For example, if you were using the Read from File source stage, your dataflow would look like this:



Candidate Finder obtains the candidate records that will form the set of potential matches that Transactional Match will evaluate later in the dataflow.

5. Double-click the Candidate Finder stage on the canvas.
6. In the **Connection** field, select the database you want to query to find candidate records. If the database you want is not listed, open Management Console and define the database connection there first.
7. In the SQL field, enter a SQL `SELECT` statement that finds records that are candidates based on the value in one of the dataflow fields. To reference dataflow fields, use the format `${FieldName}`, where `FieldName` is the name of the field you want to reference.

For example, if you wanted to find records in the database where the value in the `LastName` column is the same as the dataflow records' `Customer_LastName` field, you would write a SQL statement like this:

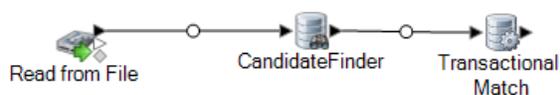
```
SELECT FirstName, LastName, Address, City, State, PostalCode
FROM Customer_Table
WHERE LastName = ${Customer_LastName};
```

- On the **Field Map** tab, select which fields in the dataflow should contain the data from each database column.

The **Selected Fields** column lists the database columns and the **Stage Fields** lists the fields in the dataflow.

- Click **OK**.
- Drag a Transactional Match stage onto the canvas and connect the Candidate Finder stage to it.

For example, if you are using a Read from File input stage your dataflow would now look like this:



Transactional Match matches suspect records against candidate records that are returned from the Candidate Finder stage. Transactional Match uses matching rules to compare the suspect record to all candidate records with the same candidate group number (assigned in Candidate Finder) to identify duplicates.

- Double-click the Transactional Match stage on the canvas.
- In the **Load match rule** field, select one of the predefined match rules which you can either use as-is or modify to suit your needs. If you want to create a new match rule without using one of the predefined match rules as a starting point, click **New**. You can only have one custom rule in a dataflow.

Note: The Dataflow Options feature in Enterprise Designer enables the match rule to be exposed for configuration at runtime.

- For information about modifying the other options, see [Building a Match Rule](#) on page 69.
- When you are done configuring the Transactional Match stage, click **OK**.
- Drag a sink stage onto the canvas and connect it to the Transactional Match stage.

For example, if you were using a Write to File sink stage your dataflow would look like this:



- Double-click the sink stage and configure it.

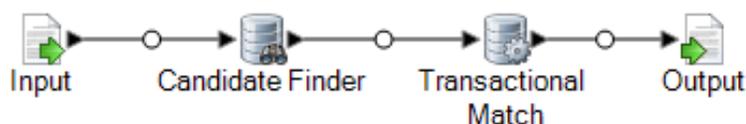
For information on configuring sink stages, see the *Dataflow Designer's Guide*.

You now have a dataflow that will match records from two data sources.

Example of Matching Records Against a Database

As a sales executive for an online sales company you want to determine if an online prospect is an existing customer or a new customer.

The following dataflow service provides a solution to the business scenario:



This dataflow is a service that evaluates prospect data sent to it by an API call or web service call. It evaluates the data against customer data in a customer database to determine if a prospect is a customer.

The Input stage is configured so that the dataflow accepts the following input fields: AddressLine1, City, Name, PostalCode, and StateProvince. AddressLine1 and Name are the fields that are key to the dataflow processing in this template.

The Candidate Finder stage obtains the candidate records that will form the set of potential matches that the Transactional Match stage will evaluate.

The Transactional Match stage matches suspect records against potential candidate records that are returned from the Candidate Finder Stage. Transactional Match uses matching rules to compare the suspect record to all candidate records with the same candidate group number (assigned in Candidate Finder) to identify duplicates. In this example, Transactional Match compares LastName and AddressLine1.

The Output stage returns the results of the dataflow through an API or web service response.

Matching Records Using Multiple Match Rules

If you have records that you want to match and you want to use more than one matching operation, you can create a dataflow that uses more than one match key then combines the results to effectively match on multiple separate criteria. For example, say you want to create a dataflow that matches records where:

The name and address match

OR

The date of birth and government ID match

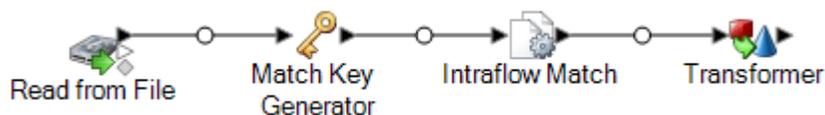
To perform matching using this logic, you create a dataflow that performs name and address matching in one stage, and date of birth and government ID matching in another stage, then combine the matching records into a single collection.

This topic provides a general procedure for setting up a dataflow where matching occurs over the course of two matching stages. For purposes of illustration this procedure uses Intraflow Match stages. However, you can use this technique with Interflow Match as well.

1. In Enterprise Designer, create a new dataflow.
2. Drag a source stage onto the canvas.
3. Double-click the source stage and configure it. See the *Dataflow Designer's Guide* for instructions on configuring source stages.
4. Define the first matching pass. The results of this first matching pass will be collections of records that match on your first set of matching criteria, for example records that match on name and address.
 - a) Drag a Match Key Generator and Intraflow Match stage to the canvas and connect them so that you have a dataflow that looks like this:

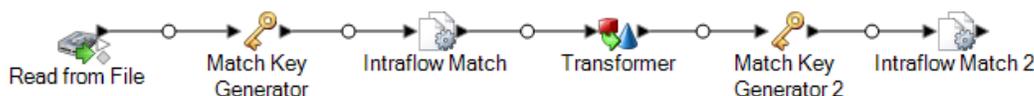


- a) In the Match Key Generator stage, define the match key to use for the first matching pass. For example, if you want the first matching pass to match on name and address, you may create a match key based on the fields containing the last name and postal code.
 - b) In the Intraflow Match stage, define the match rules you want to perform the first matching pass. For example, if you may configure this matching stage to match on name and address.
5. Save the collection numbers from the first matching pass to another field. This is necessary because the CollectionNumber field will be overwritten during the second matching pass. It is necessary to rename the CollectionNumber field in order to preserve the results of the first matching pass.
 - a) Drag a Transformer stage to the canvas and connect it to the Intraflow Match stage so that you have a dataflow that looks like this:



- b) Configure the Transformer stage to rename the field CollectionNumber to CollectionNumberPass1.
6. Define the second matching pass. The results of this second matching pass will be collections of records that match on your second set of matching criteria, for example records that date of birth and government ID.

- a) Drag a Match Key Generator and Intraflow Match stage to the canvas and connect them so that you have a dataflow that looks like this:



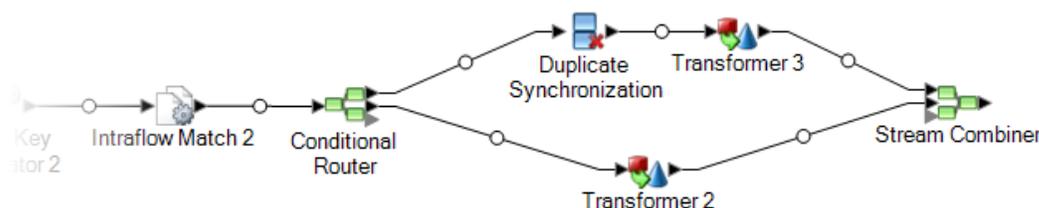
- b) In the second Match Key Generator stage, define the match key to use for the second matching pass.

For example, if you want the second matching pass to match date of birth and government ID, you might create a match key based on the fields containing the birthday and government ID.

- c) In the second Intraflow Match stage, define the match rule for the second matching pass.
For example, if you may configure this matching stage to match on date of birth and government ID.

7. Determine if any of the duplicate records identified by the second matching pass were also identified as duplicates in the first matching pass.

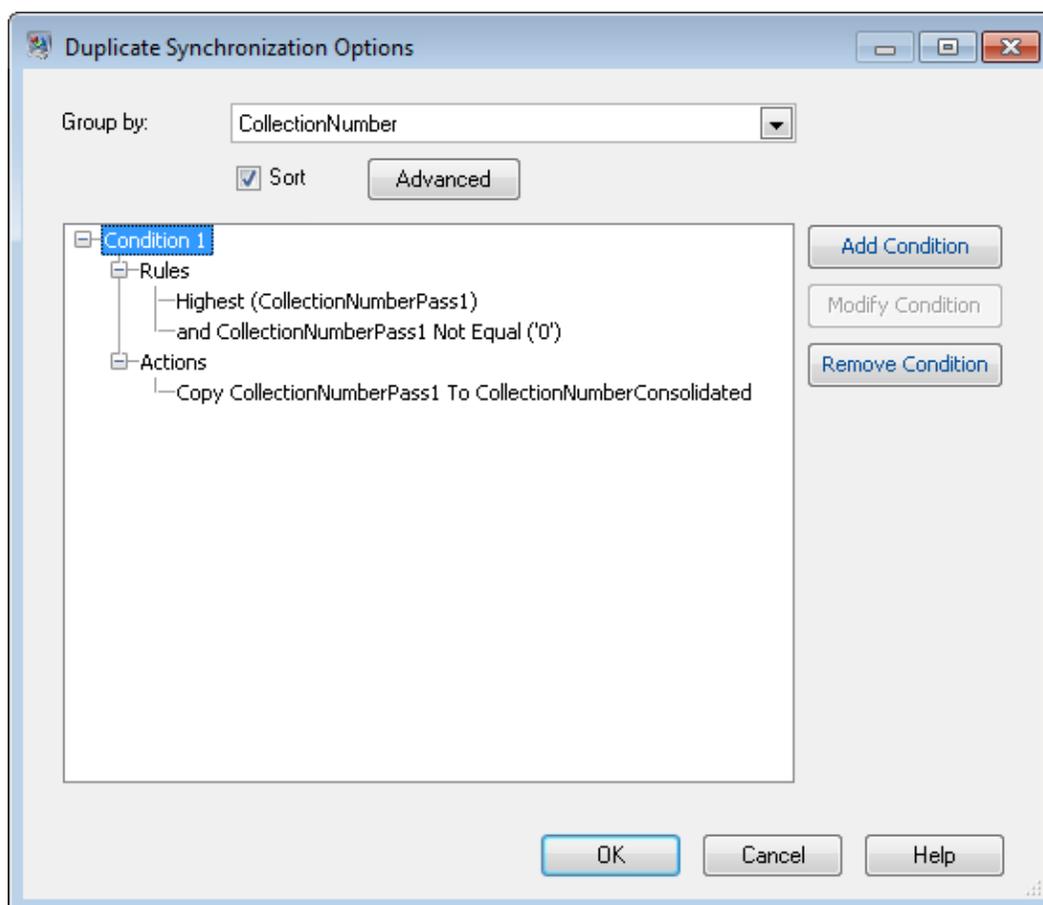
- a) Create the dataflow snippet shown below following the second Intraflow Match stage:



- b) Configure the Conditional Router stage so that records where the CollectionNumber field is not equal to 0 are routed to the Duplicate Synchronization stage.

This will route the duplicates from the second matching pass to the Duplicate Synchronization stage.

- c) Configure the Duplicate Synchronization stage to group records by the CollectionNumber field (this is the collection number from the second matching pass). Then within each collection, identify whether any of the records in the collection were also identified as duplicates in the first matching pass. If they were, copy the collection number from the first pass to a new field called CollectionNumberConsolidated. To accomplish this, configure Duplicate Synchronization as shown here:



- d) In the Transformer stage that follows the Duplicate Synchronization stage, create a custom transform using this script:

```
if (data['CollectionNumberConsolidated'] == null) {
  data['CollectionNumberConsolidated'] = data['CollectionNumber']
}
```

- e) In the Transformer that immediately follows the Conditional Router (Transformer 2 in sample dataflow) configure a transform to copy CollectionNumberPass1 to CollectionNumberConsolidated.

This takes the unique records from the second matching pass and copies CollectionNumberPass1 to CollectionNumberConsolidated.

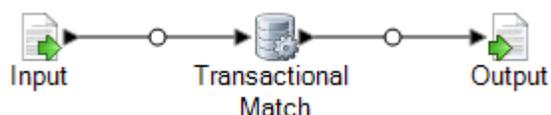
8. After the Stream Combiner you will have collections of records that match in either of the matching passes. The CollectionNumberConsolidated field indicates the matching records. You can add a sink or any additional processing you wish to perform after the Stream Combiner stage.

Creating a Universal Matching Service

A universal matching service is a service that can use any of your match rules to perform matching and can accept any input fields. The service takes a match rule name as an input option, allowing you specify the match rule you want to use in the API call or web service request. The service does not have a predefined input schema so you can include whatever fields are appropriate for the type of records you want to match. By creating a universal matching service you can avoid having separate services for each match rule, enabling you to add new match rules without having to add a service.

This procedure shows how to create a universal matching service and includes an example of a web service request to the universal matching service.

1. In Enterprise Designer, create a new service dataflow.
2. Drag an Input stage, a Transactional Match stage, and an Output stage to the canvas and connect them so that you have a dataflow that looks like this:



3. Double-click the Transactional Match stage.
4. In the **Load match rule** field, select any match rule. For example, you can select the default **Household** match rule.

Even though you will specify the match rule in the service request, you have to configure the Transactional Match stage with a default match rule in order for the dataflow to be valid. If you do not select a match rule the dataflow will fail validation and you will not be able to expose it.

5. Click **OK**.
6. Double-click the Output stage.
7. Choose to expose the fields `MatchRecordType` and `MatchScore`.
8. Click **OK**.

Note: There is no need to expose any fields in the Input stage since input fields will be specified as user-defined fields in the service request.

9. Click **Edit > Dataflow Options**.
10. Click **Add**.
11. Expand **Transactional Match** and check the box next to **Match Rule**.

This exposes the match rule option as a run-time option, making it possible to specify the match rule in the service request.

12. Click **OK** then click **OK** again to close the **Dataflow Options** window.

13. Save and expose the dataflow.

You now have a universal match service that you can use to perform matching using any of the match rules defined in the Match Rules Management tool in Enterprise Designer. When calling the service, specify the match rule in the `MatchRule` option and specify the input fields as user-defined fields.

Example: Calling the Universal Matching Service

You have created a match rule named `AddressAndBirthday` in the Match Rules Management tool. This match rule matches records using the fields `Address` and `Birthday`. You want to use the universal matching service to perform matching using this rule through a SOAP web service request.

To accomplish this, you would have a SOAP request that specifies `AddressAndBirthday` in the `MatchRule` element and the record's fields in the `user_fields` element.

```
<soapenv:Envelope xmlns:soapenv="http://schemas.xmlsoap.org/soap/envelope/"
xmlns:univ="http://www.precisely.com/spectrum/services/UniversalMatchingService">
```

```
  <soapenv:Header/>
  <soapenv:Body>
    <univ:UniversalMatchingServiceRequest>
      <univ:options>
        <univ:MatchRule>AddressAndBirthday</univ:MatchRule>
      </univ:options>
      <univ:Input>
        <univ:Row>
          <univ:user_fields>
            <univ:user_field>
              <univ:name>Name</univ:name>
              <univ:value>Bob Smith</univ:value>
            </univ:user_field>
            <univ:user_field>
              <univ:name>Address</univ:name>
              <univ:value>4200 Parliament Pl</univ:value>
            </univ:user_field>
            <univ:user_field>
              <univ:name>Birthday</univ:name>
              <univ:value>1973-6-15</univ:value>
            </univ:user_field>
          </univ:user_fields>
        </univ:Row>
        <univ:Row>
          <univ:user_fields>
            <univ:user_field>
              <univ:name>Name</univ:name>
              <univ:value>Robert M. Smith</univ:value>
            </univ:user_field>
            <univ:user_field>
              <univ:name>Address</univ:name>
```

```

        <univ:value>4200 Parliament Pl</univ:value>
      </univ:user_field>
      <univ:user_field>
        <univ:name>Birthday</univ:name>
        <univ:value>1973-6-15</univ:value>
      </univ:user_field>
    </univ:user_fields>
  </univ:Row>
</univ:Row>
  <univ:user_fields>
    <univ:user_field>
      <univ:name>Name</univ:name>
      <univ:value>Bob Smith</univ:value>
    </univ:user_field>
    <univ:user_field>
      <univ:name>Address</univ:name>
      <univ:value>424 Washington Blvd</univ:value>
    </univ:user_field>
    <univ:user_field>
      <univ:name>Birthday</univ:name>
      <univ:value>1959-2-19</univ:value>
    </univ:user_field>
  </univ:user_fields>
</univ:Row>
</univ:Input>
</univ:UniversalMatchingServiceRequest>
</soapenv:Body>
</soapenv:Envelope>

```

This request would result in the following response:

```

<soap:Envelope xmlns:soap="http://schemas.xmlsoap.org/soap/envelope/">
  <soap:Body>
    <ns3:UniversalMatchingServiceResponse
xmlns:ns2="http://spectrum.precisely.com/"
xmlns:ns3="http://www.precisely.com/spectrum/services/UniversalMatchingService">
      <ns3:Output>
        <ns3:Row>
          <ns3:MatchScore/>
          <ns3:MatchRecordType>Suspect</ns3:MatchRecordType>
          <ns3:user_fields>
            <ns3:user_field>
              <ns3:name>Name</ns3:name>
              <ns3:value>Bob Smith</ns3:value>
            </ns3:user_field>
            <ns3:user_field>
              <ns3:name>Birthday</ns3:name>
              <ns3:value>1973-6-15</ns3:value>
            </ns3:user_field>
          </ns3:user_fields>
        </ns3:Row>
      </ns3:Output>
    </ns3:UniversalMatchingServiceResponse>
  </soap:Body>
</soap:Envelope>

```

```

        <ns3:user_field>
          <ns3:name>Address</ns3:name>
          <ns3:value>4200 Parliament Pl</ns3:value>
        </ns3:user_field>
      </ns3:user_fields>
    </ns3:Row>
    <ns3:Row>
      <ns3:MatchScore>100</ns3:MatchScore>
      <ns3:MatchRecordType>Duplicate</ns3:MatchRecordType>
      <ns3:user_fields>
        <ns3:user_field>
          <ns3:name>Name</ns3:name>
          <ns3:value>Robert M. Smith</ns3:value>
        </ns3:user_field>
        <ns3:user_field>
          <ns3:name>Birthday</ns3:name>
          <ns3:value>1973-6-15</ns3:value>
        </ns3:user_field>
        <ns3:user_field>
          <ns3:name>Address</ns3:name>
          <ns3:value>4200 Parliament Pl</ns3:value>
        </ns3:user_field>
      </ns3:user_fields>
    </ns3:Row>
  </ns3:Output>
</ns3:UniversalMatchingServiceResponse>
</soap:Body>
</soap:Envelope>

```

Using an Express Match Key

Express key matching can be a useful tool for reducing the number of compares performed and thereby improving execution speed in dataflows that use an Interflow Match or Intraflow Match stage. If two records have an exact match on the express key, the candidate is considered a 100% match and no further matching attempts are made. If two records do not match on an express key value, they are compared using the rules-based method. However, a loose express key results in many false positive matches.

1. Open your dataflow in Enterprise Designer.
2. Double-click the Match Key Generator stage.
3. Check the box **Generate express match key**.
4. Click **Add**.
5. Complete the following fields:

Table 6: Match Key Generator Options

Option Name	Description and Valid Values
-------------	------------------------------

Algorithm	
-----------	--

Option Name	Description and Valid Values
-------------	------------------------------

Specifies one of these algorithms to use to generate the match key:

Consonant	Returns specified fields with consonants removed.
Double Metaphone	Returns a code based on a phonetic representation of their characters. Double Metaphone is an improved version of the Metaphone algorithm, and attempts to account for the many irregularities found in different languages.
Koeln	Indexes names by sound as they are pronounced in German. Allows names with the same pronunciation to be encoded to the same representation so that they can be matched, despite minor differences in spelling. The result is always a sequence of numbers; special characters and white spaces are ignored. This option was developed to respond to limitations of Soundex.
MD5	A message digest algorithm that produces a 128-bit hash value. This algorithm is commonly used to check data integrity.
Metaphone	Returns a Metaphone coded key of selected fields. Metaphone is an algorithm for coding words using their English pronunciation.
SpanishMetaphone	Returns a Metaphone coded key of selected fields for the Spanish language. This metaphone algorithm codes words using their Spanish pronunciation.
Metaphone 3	Improves upon the Metaphone and Double Metaphone algorithms with more exact consonant and internal vowel settings that allow you to produce words or names more or less closely matched to search terms on a phonetic basis. Metaphone 3 increases the accuracy of phonetic encoding to 98%. This option was developed to respond to limitations of Soundex.
Nysis	Phonetic code algorithm that matches an approximate pronunciation to an exact spelling and indexes words that are pronounced similarly. Part of the New York State Identification and Intelligence System. Say, for example, that you are looking for someone's information in a database of people. You believe that the person's name sounds like "John Smith", but it is in fact spelled "Jon Smyth". If you conducted a search looking for an exact match for "John Smith" no results would be returned. However, if you index the database using the NYSIIS algorithm and search using the NYSIIS algorithm again, the correct match will be returned because both "John Smith" and "Jon Smyth" are indexed as "JAN SNATH" by the algorithm.
Phonix	Preprocesses name strings by applying more than 100 transformation rules to single characters or to sequences of several characters. 19 of those rules are applied only if the characters are at the beginning of the string, while 12 of the rules are applied only if they are at the middle of the string, and 28 of the rules are applied only if they are at the end of

Option Name	Description and Valid Values
	<p>the string. The transformed name string is encoded into a code that is comprised by a starting letter followed by three digits (removing zeros and duplicate numbers). This option was developed to respond to limitations of Soundex; it is more complex and therefore slower than Soundex.</p>
	<p>Sonnex This algorithm determines the similarity between two French-language strings based on the phonetic representation of their characters. It returns a Sonnex coded key of the selected fields.</p>
	<p>Soundex Returns a Soundex code of selected fields. Soundex produces a fixed-length code based on the English pronunciation of a word.</p>
	<p>Substring Returns a specified portion of the selected field.</p>
Field name	<p>Specifies the field to which you want to apply the selected algorithm to generate the match key. For example, if you select a field called LastName and you choose the Soundex algorithm, the Soundex algorithm would be applied to the data in the LastName field to produce a match key.</p>
Start position	<p>Specifies the starting position within the specified field. Not all algorithms allow you to specify a start position.</p>
Length	<p>Specifies the length of characters to include from the starting position. Not all algorithms allow you to specify a length.</p>
Remove noise characters	<p>Removes all non-numeric and non-alpha characters such as hyphens, white space, and other special characters from an input field.</p>
Sort input	<p>Sorts all characters in an input field or all terms in an input field in alphabetical order.</p> <p>Characters Sorts the characters values from an input field prior to creating a unique ID.</p> <p>Terms Sorts each term value from an input field prior to creating a unique ID.</p>

6. Click **OK**.

7. If you want to specify an additional field and/or algorithm to use in generating an express match key, click **Add**, otherwise click **OK**.
8. Double-click the Interflow Match or Intraflow Match stage on the canvas.
9. Select the option **Express match on** and choose the field **ExpressMatchKey**.
This field contains the express match key produced by Match Key Generator.
10. Click **OK**.
11. Save and run your dataflow.

To determine whether a candidate was matched using an express key, look at the value of the **ExpressKeyIdentified** field, which is either Y for a match or N for no match. Note that suspect records always have an **ExpressKeyIdentified** value of N.

Analyzing Match Results

The Match Analysis tool in Enterprise Designer displays the results of one or more matching stages of the same type. The tool provides summary matching results for a dataflow and also allows you to view matching results on a record-by-record basis. You can use this information to troubleshoot or fine-tune your match rules to produce the results you want.

The Match Analysis tool provides the following features:

- Match Summary Results: Displays summary record counts for a single match result or comparisons between two match results.
- Lift/Drop charts: Uses bar charts to display an increase or decrease in matches..
- Match rules: Displays the match rules used for a single match result or the changes made to the match rules when comparing two match results.
- Match Detail results: Displays record processing details for a single match result or the comparison between two match results.

Viewing a Summary of Match Results

The Match Analysis tool can display summary of the matching processes in a dataflow, such as the number of duplicate records and the average match score. You can view the results of a single job or you can compare results between multiple jobs.

1. In **Enterprise Designer**, open the dataflow you want to analyze.
2. For each Interflow Match, Intraflow Match, or Transactional match stage whose matching you want to analyze, double-click the stage and select the **Generate data for analysis** check box.

Important: Enabling the **Generate data for analysis** option reduces performance. You should turn this option off when you are finished using the Match Analysis tool.

3. Select **Run > Run Current Flow**

Note: For optimal results, use data that will produce 100,000 or fewer records. The more match results, the slower the performance of the Match Analysis tool.

4. When the dataflow finishes running, select **Tools > Match Analysis**.

The **Browse Match Results** dialog box displays with a list of dataflows that have match results that be viewed in the Match Analysis tool. If the job you want to analyze is not listed, open the dataflow and make sure that the matching stage has the **Generate data for analysis** check box selected.

Tip: If there are a large number of dataflows and you want to filter the dataflows, select a filter option from the **Show only jobs where** drop-down list.

5. Click the "+" icon next to the dataflow you want to view to expand it.

6. Under the dataflow there is one entry for each matcher stage in the dataflow. Select the stage whose results you want to view and click **Add**.

The Match Analysis tool appears at the bottom of the Enterprise Designer window.

7. If you want to compare the matcher results side by side with the results from another matcher:

- a) Click **Add**.
- b) Select the matcher whose results you want to compare.
- c) Click **Add**.
- d) In the dataflow list, select the matcher you just added and click **Compare**.

The **Summary** tab lists matching statistics for the job. Depending on the type of matching stage used in the dataflow, you will see different information.

For Intraflow Match you will see the following summary information:

Input Records	The total number of records processed by the matcher stage.
Duplicate Records	Number of records that match another record within a match group.
Unique Records	A suspect or candidate record that does not match any other records in a match group. If it is the only record in a match group, a suspect is automatically unique.
Match Groups	(Group By) Records grouped together either by a match key or a sliding window.
Duplicate Collections	A duplicate collection consists of a Suspect and its Duplicate records grouped together by a CollectionNumber. Unique records always belong to CollectionNumber 0.
Express Matches	An express match is made when a suspect and candidate have an exact match on the contents of a designated field, usually an ExpressMatchKey provided by the Match Key Generator. If an Express Match is made no

further processing is done to determine if the suspect and candidate are duplicates.

Average Score The average match score of all duplicates. The possible values are 0-100, with 0 indicating a poor match and 100 indicating an exact match.

For Interflow Match you will see the following summary information:

Duplicate Collections A duplicate collection consists of a Suspect and its Duplicate records grouped together by a CollectionNumber. Unique records always belong to CollectionNumber 0.

Express Matches An express match is made when a suspect and candidate have an exact match on the contents of a designated field, usually an ExpressMatchKey provided by the Match Key Generator. If an Express Match is made no further processing is done to determine if the suspect and candidate are duplicates.

Average Score The average match score of all duplicates. The possible values are 0-100, with 0 indicating a poor match and 100 indicating an exact match.

Input Suspects The number of records in the input stream that the matcher tried to match to other records.

Suspects with Duplicates The number of input suspects that matched at least one candidate record.

Unique Suspects The number of input suspects that did not match any candidate records.

Suspects with Candidates The number of input suspects that had at least one candidate record in its match group and therefore had at least one match attempt.

Suspects without Candidates The number of input suspects that had no candidate records in its match group and therefore had no match attempts.

For Transactional Match, you will see the following summary information:

Average Score The average match score of all duplicates. The possible values are 0-100, with 0 indicating a poor match and 100 indicating an exact match.

Input Suspects The number of records in the input stream that the matcher tried to match to other records.

Suspects with Duplicates The number of input suspects that matched at least one candidate record.

Unique Suspects The number of input suspects that did not match any candidate records.

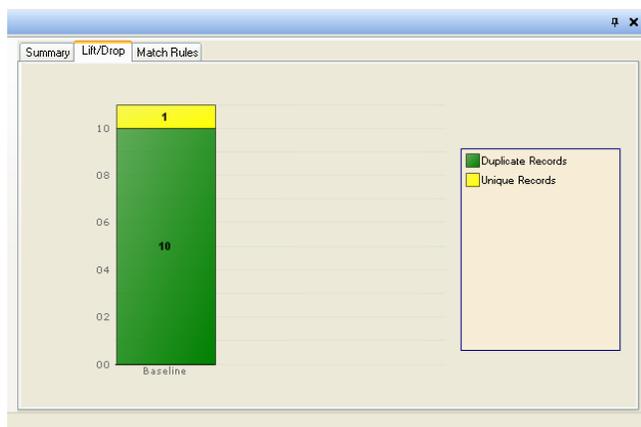
Suspects with Candidates The number of input suspects that had at least one candidate record in its match group and therefore had at least one match attempt.

Suspects without Candidates The number of input suspects that had no candidate records in its match group and therefore had no match attempts.

The **Lift/Drop** tab of the Match Analysis tool displays duplicate and unique record counts in a bar chart for the selected baseline and, optionally, comparison results. *Lift* is the increase in the number

of duplicate records. *Drop* is the decrease in the number of duplicate records. Unique records are shown in yellow and duplicate records are shown in green.

If only a baseline job is selected, the chart will show the results for that one job:



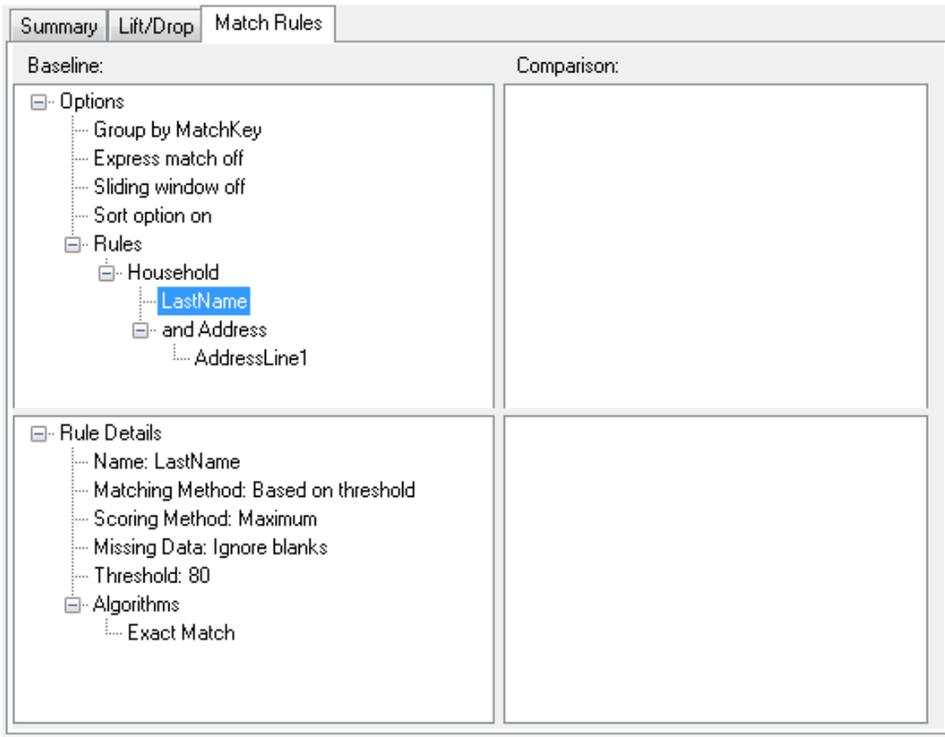
If both a baseline and a comparison job are selected, a chart for the baseline and comparison jobs are shown side by side:



The **Match Rules** tab of the Match Analysis tool displays the match rules used for a single match result or the changes made to the match rules when comparing two match results.

Match rules are displayed in a hierarchical structure similar to how they are displayed in the stage in which they were created. The rule hierarchy contains two nodes: Options and Rules. The Options node shows the stage settings for the selected match result. The Rules node shows the match rules for the selected match result.

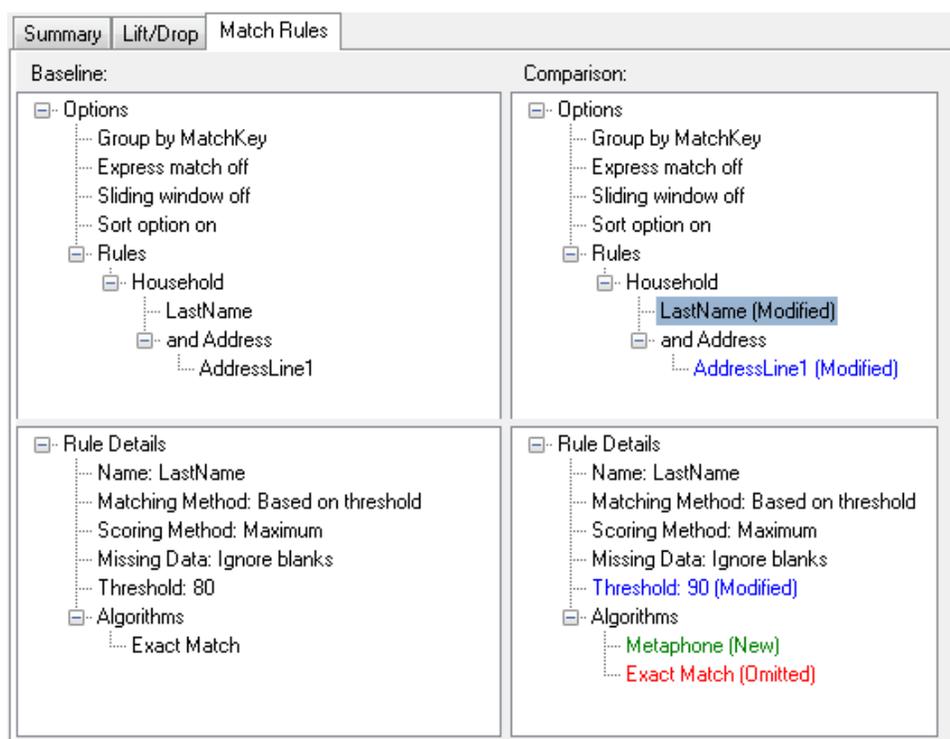
To view rule details, select a node in the hierarchy.



If you are comparing match rules between multiple jobs, differences between the baseline and comparison match results are color coded as follows:

- Blue** Indicates that the match rule in the comparison match result was modified.
- Green** Indicates that the match rule in the comparison match result was added.
- Red** Indicates that the match rule in the comparison match result was omitted.

This is demonstrated in this figure.



Viewing Record-Level Match Results

Detailed results displays a collection of details about match records for match results set.

To display detailed results:

1. In the Match Analysis tool, specify a baseline job and, optionally, a comparison job.
2. Click **Details**.

The baseline match results are displayed based on the selected view in the Show drop-down list. The following table lists the columns displayed for each match stage type.

Table 7: Detailed Results Data Displayed

Detail-Related Results	Intraflow	Interflow	Transactional
Input Record Number	X	X	X
Match Group	X	X	

Detail-Related Results	Intraflow	Interflow	Transactional
Express Key	X	X	
Express Key Driver Record	X	X	
Collection Number	X	X	X
Match Record Type	X	X	X
Fields used by the rules	X	X	X
Overall (top-level) rule score	X		
Candidate Group		X	X
Match Score Select a match results in the Match Results List and then click Remove .		X	X

For information about the match rate chart, see [Match Rate Chart](#) on page 119.

- In the **Analyze** field, choose one of the following:
 - Baseline** Displays the match results from the baseline run.
 - Comparison** Displays the match results of the comparison run.
- Select one of the following values from the **show** list and then click **Refresh**. If you are analyzing baseline results, the options are:
 - Suspects with Candidates: (All matchers) Displays suspect records and all candidate records that attempted to match to each suspect.
 - Suspects with Duplicates: (All matchers) Displays all suspect records and candidate records that matched to each suspect.
 - Suspects with Express Matches: (Interflow Match and Intraflow Match, when Express Match Key is enabled) Displays suspect and candidate records that match based on the Express Match Key.

- Duplicate Collections: (Intraflow and Interflow) Displays all duplicate collections by collection number.
- Match Groups: (Intraflow and Interflow) Displays records by match groups.
- Candidate Groups: (Transactional Match) Displays records by candidate groups.
- Unique Suspects: (Interflow and Transactional Match) Displays all suspect records that did not match to any candidate records.
- Unique Records: (Intraflow) Displays all non-matched records.
- Suspects without Candidates: (Interflow and Transactional Match) Displays all suspects that contained no candidates to match against.
- All Records: Displays all records processed by the matching stage.

If you are analyzing comparison results, the show options are:

- New Matches: (Intraflow) Displays all new matches and its related suspects. This view combines the results of Suspects with New Duplicates and New Suspects into one view.
- New Matched Suspects: (Interflow and Transactional Match) Displays suspects that had no duplicates in the baseline but have at least one duplicate in the comparison.
- New Unique Suspects: (Interflow and Transactional Match) Displays suspects that had duplicates in the baseline but have none in the comparison.
- Missed Matches: (Intraflow) Displays all missed matches. This view combines the results of Suspects with Missed Duplicates and Missed Suspects into one view.
- Suspects with New Duplicates: (All matchers) Displays records that are new duplicates for records that were suspects in the baseline and remained suspects in the comparison.
- Suspects with Missed Duplicates: (All matchers) Displays records that are missed duplicates for records that were suspects in the baseline and remained suspects in the comparison.
- New Suspects: (Intraflow) Displays records that are suspects in the comparison match result, but were not Suspects in the baseline.
- Missed Suspects (Intraflow) Displays records that are not suspects in the comparison result, but were suspects in the baseline.

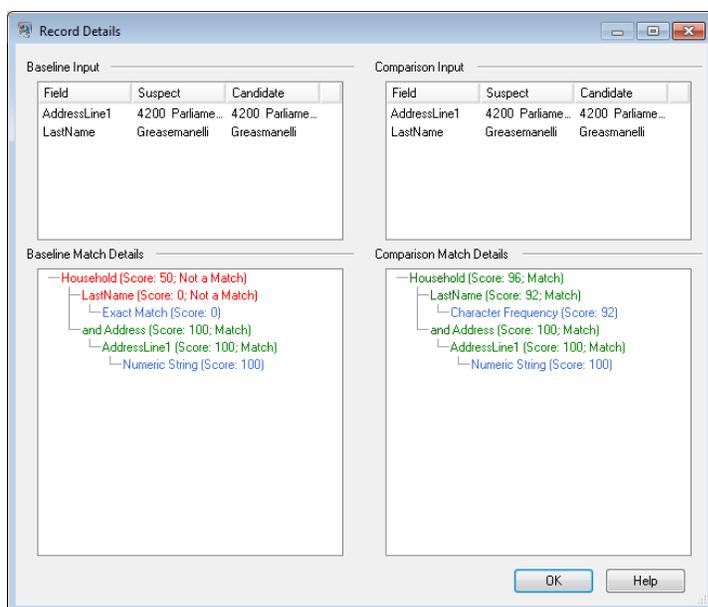
5. Expand a suspect record to view its candidates.
6. Select a candidate record and click **Details**.

Note: This option is not available when Sliding Window is enabled in Intraflow Match stages.

The **Record Details** dialog box shows field-level data as well as the record's match score for each match rule. If you specified both a baseline and a comparison job run, you can see the record's results for both baseline and comparison runs.

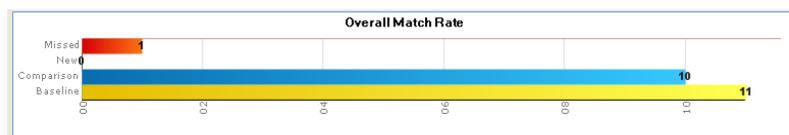
- Baseline Input—Displays the field level data, from both the suspect and candidate, used in the match.
- Baseline Match Details—Displays scoring information for each node in the match rules.
- Comparison Input—Displays the field level data, from both the suspect and candidate, used in the match.

- **Comparison Match Details**—Displays scoring information for each node in the match rules. Green text represents a match for a node in the rules. Red text represents a non-match for a node in the rules.



Match Rate Chart

Match Rate charts graphically display match information in detail views.



For Intraflow matches, it displays one chart displaying overall matches:

- **Baseline Matches:** Total number of matches in the baseline result.
- **Comparison Matches:** Total number of matches in the comparison result.
- **New Matches:** A count of all records that were unique in the baseline result, but are a suspect or duplicate in the comparison result.
- **Missed Matches:** A count of all records that were suspects or duplicates in the baseline result, but are unique in the comparison result.

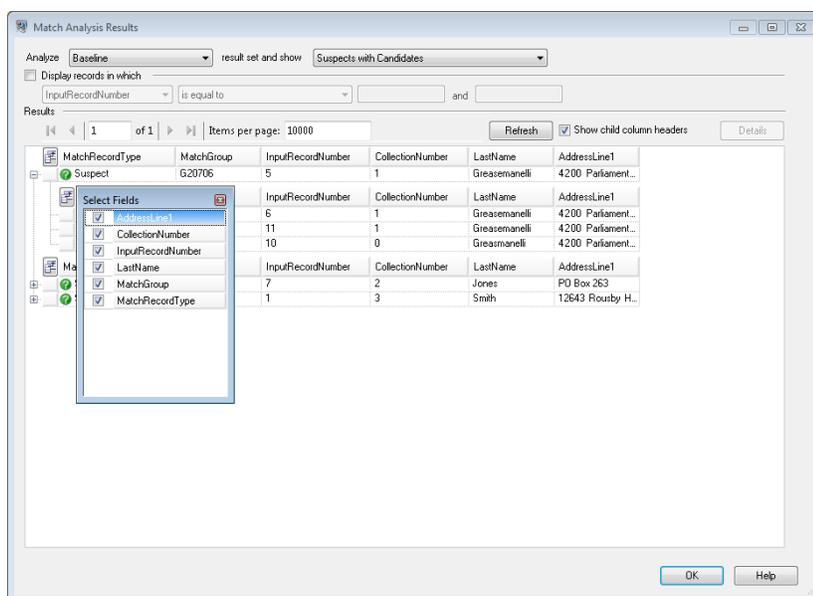
For Interflow and Transactional matches, it displays two charts:

- **Overall Match Rate**
- **Baseline Matches:** Total number of matches in the baseline result.
- **Comparison Matches:** Total number of matches in the comparison result.
- **New Matches:** A count of all records that were unique in the baseline result, but are a suspect or duplicate in the comparison result.

- Missed Matches: A count of all records that were suspects or duplicates in the baseline result, but are unique in the comparison result.
- Suspect Match Rate
- Baseline Matches: A count of all Suspects that were not unique in the baseline.
- Comparison Matches: A count of all suspects that were not unique in the comparison.
- New Matches: A count of all suspects that were unique in the baseline, but are matches in the comparison result.
- Missed Matches: A count of all suspects that were matches in the baseline, but are unique in the comparison result.

Using Field Chooser

Click the Field Chooser icon  to display selected columns in the Match Analysis Results. Field Chooser displays at the parent level and the child level. You can independently select display columns for parents and children.



Filtering Records

Use the **Display records in which** check box to filter the detail match records displayed. You can filter records based on several operators to compare user-provided values against data in one field of each detail match record.

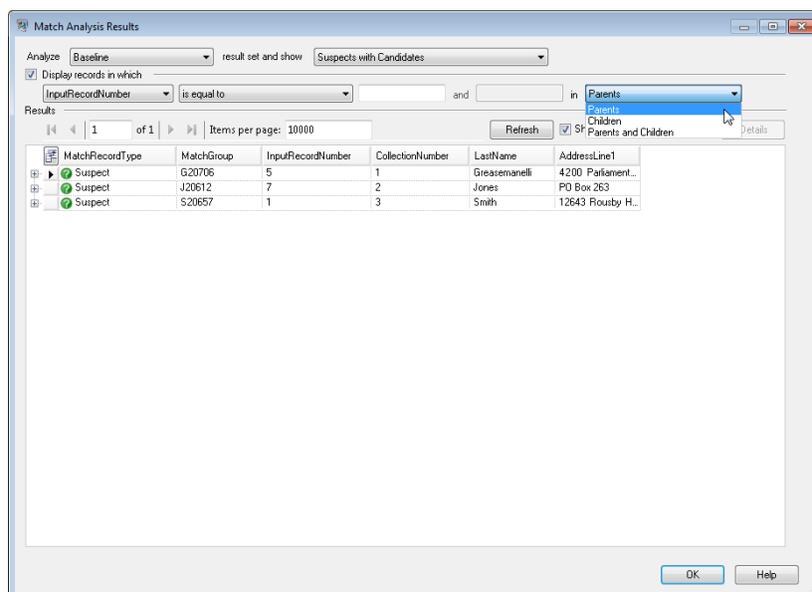
The operators you can choose are:

- String-type fields (MatchGroup, MatchRecordType, any matching data)
- contains
- is between
- is equal to

- is not equal to
- starts with
- Numeric-type fields (CollectionNumber, InputRecordNumber, MatchScore)
- is between
- is equal to
- is not equal to
- is greater than
- is greater than or equal to
- is less than
- is less than or equal to

To filter records:

1. Select a baseline or comparison match result from the Match Analysis Results view and click **Refresh**.
2. Select the **Display records in which** check box.



3. Select a field from the Field list box.
4. Select an operator.
5. Type a value for the selected operator type. If you select **is between**, type a range of values.
6. When filtering on suspect views, you can filter on:
 - Parents—Filter just on parents (Suspects), all children returned.
 - Children—Filter out any children that do not fall in the filter range. Parent (Suspect) nodes returned.
 - Parents and Children—Filter on parents (Suspects), then if any parents are returned, filter on its children

7. Click **Refresh**. Records that fall in the range of the options and values are displayed. If no records fall in the range of the selected options and values, a message displays that no records were returned.

Analyzing Match Rule Changes

You can use the Match Analysis tool in Enterprise Designer to view in detail the effect that a change in a match rule has in the dataflow match results. You can do this by running the dataflow, making changes, re-running the dataflow, and then viewing the results in the Match Analysis tool. This procedure describes how to do this.

Important: When comparing match results, the input data used for the baseline and comparison runs must be identical. Using different input data can cause misleading results. Observe the following to help ensure an accurate comparison:

- Use the same input files or tables
- Sort the data in the same way prior to the matching stage
- Use the same Candidate Finder queries when using Transactional Match

1. In **Enterprise Designer**, open the dataflow you want to analyze.
2. For each Interflow Match, Intraflow Match, or Transactional match stage whose matching you want to analyze, double-click the stage and select the **Generate data for analysis** check box.

Important: Enabling the **Generate data for analysis** option reduces performance. You should turn this option off when you are finished using the Match Analysis tool.

3. Select **Run > Run Current Flow**

Note: For optimal results, use data that will produce 100,000 or fewer records. The more match results, the slower the performance of the Match Analysis tool.

4. In the dataflows matcher stage or stages, make the match rule changes you want then run the dataflow again.

For example, if you want to test the effect of increasing the threshold value, change the threshold value and run the dataflow again.

5. When the dataflow finishes running, select **Tools > Match Analysis**.

The **Browse Match Results** dialog box displays with a list of dataflows that have match results that be viewed in the Match Analysis tool. If the job you want to analyze is not listed, open the dataflow and make sure that the matching stage has the **Generate data for analysis** check box selected.

Tip: If there are a large number of dataflows and you want to filter the dataflows, select a filter option from the **Show only jobs where** drop-down list.

- On the left side of the Match Analysis pane, there is a list of the matcher stages, one per run. Select the matcher stage in the run that you want to use as the baseline for comparison then click **Baseline**. Then, select the run you want to compare the baseline to and click **Compare**.

You can now compare summary match results, such as the total number of duplicate records, as well as detailed record-level information that shows how each record was evaluated against the match rules.

Example of Match Results Comparison

For example, say you run a job named HouseholdRelationshipsAnalysis. You want to test the effect of a change to the Household Match 2 stage. Your first run the job using the original settings, then you modify the match rules in the Household Match 2 stage and run the job again. In the Match Analysis tool, the run with a job ID of 10 is the run with the original settings, so you set it as the baseline. The run with a job ID of 13 is run with the modified match rule. When you click **Compare**, you can see that the modified match rule (job ID 13) produced one more duplicate record and one less unique record than the original match rule.

	Baseline	Comparison	Changes
Input Records	11	11	0
Duplicate Records	10	11	1
Unique Records	1	0	-1
Match Groups	3	3	0
Duplicate Collections	3	3	0
Express Matches	0	0	0
Average Score	99.8	99.8	0.0

Adding Match Results

If you run a job while the Match Analysis Tool is open and the Match Results List is empty, the match results are automatically added to the list. After a match result has been added, the Match Analysis Tool only adds match results of the same match type (Interflow Match, Intraflow Match, or Transactional Match).

If you want to analyze match results of a different type than what is currently selected in the Match Analysis Tool, follow these steps.

- Select all match results in the Match Results List and then click **Remove**.
- Open a job from the Server Explorer that uses a different matching stage or click the tab above the canvas if the job is already open.
- Run the job.

When the job finishes running, the match results from the last job instance are added to the Match Results List.

Removing Match Results

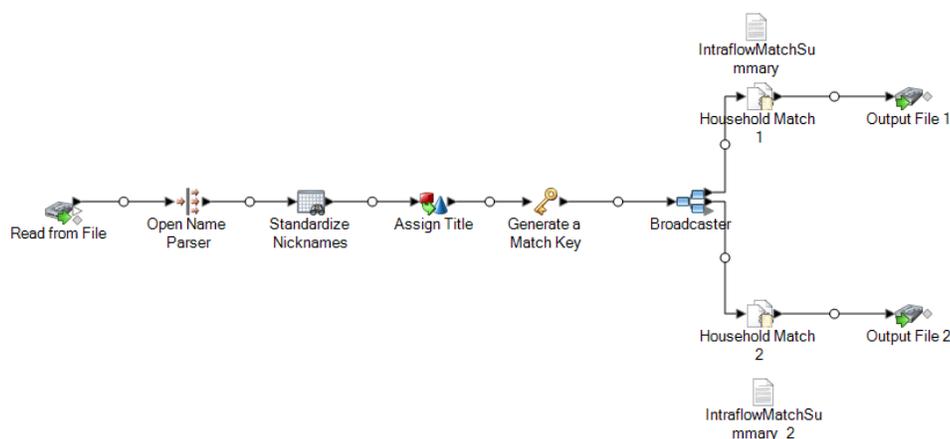
To remove a match results from the Match Results List, select a match results in the Match Results List and then click **Remove**.

The system updates the Match Results list and Summary tab as follows:

- If the removed match results was neither the Baseline nor the Comparison match results, the match results is removed and no changes to the Summary tab occur.
- If the removed match results was set as the Baseline, the system sets the next oldest match results as the new Baseline and updates the Summary tab to display the new Baseline data only.
- If the removed match results was set as the Comparison match results, the system updates the Summary tab to display the existing Baseline data only.
- If the removed match results is one of two displayed in the Match Results list, the remaining match results is set as the new Baseline and system updates the Summary tab to display the new Baseline data only.

Example: Using Match Analysis

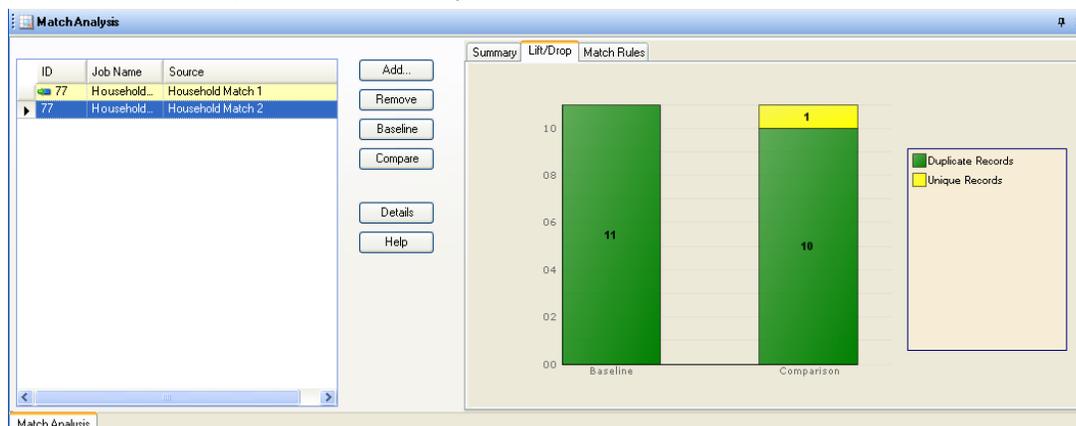
This example demonstrates how to use the Match Analysis tool to compare the lift/drop rates of two different matches. Before the data is sent through a matcher, it is split into two streams using a Broadcaster. Each stream is then sent through an Intraflow Match stage. Each data stream includes identical copies of the processed data. Each Intraflow Match stage uses different matching algorithm and generates Match Analysis data that you can use to compare the lift/drop of various matches.



This example dataflow is available in Enterprise Designer. Go to **File > New > Dataflow > From template** and select **HouseholdRelationshipsAnalysis**. This dataflow requires the following products: Advanced Matching, Data Normalization, and Universal Name. It also requires you to load the Table Lookup core database and the Open Parser base tables.

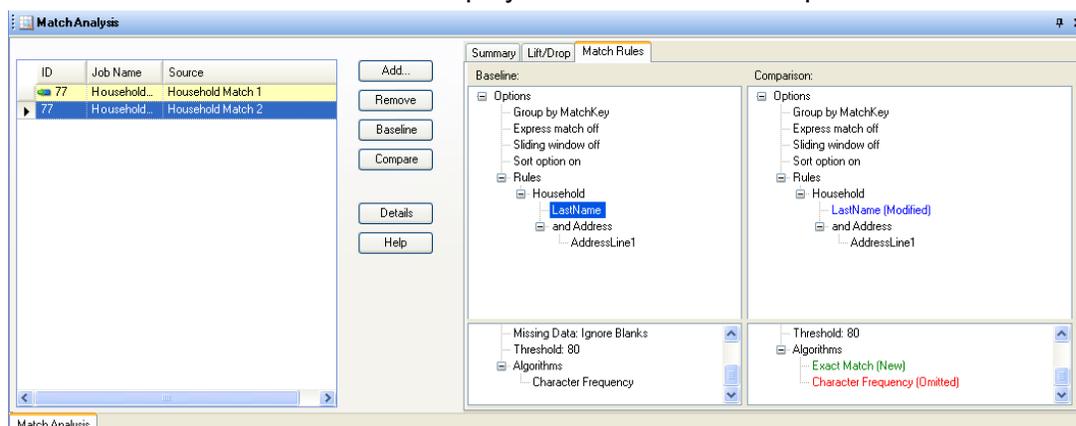
To use view this example:

1. Run the dataflow.
2. Select **Tools > Match Analysis**.
3. From **Browse Match Results** window, expand **HouseholdRelationshipAnalysis**, select **Household Match 1** and **Household Match 2** from the Source list, and then click **Add**.
4. Select **Household Match 1** in the Match Results List and click **Compare**. This displays results on the **Summary** tab.
5. Click the **Lift/Drop** tab. This displays the Lift/Drop chart.



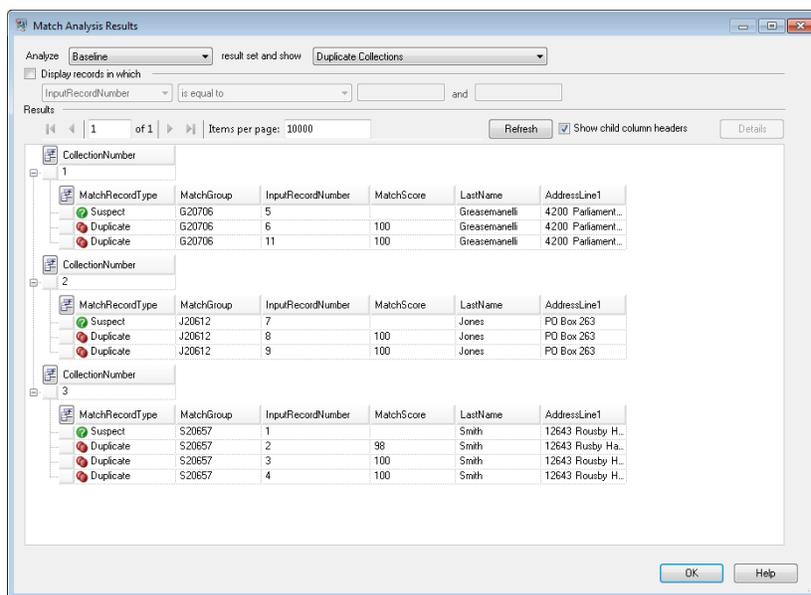
This chart shows the differences between the duplicate and unique records generated for the different match rules used.

6. Click the **Match Rules** tab. This displays the match rules comparison.



From this tab you can see that the algorithm has been changed; Character Frequency is omitted and Exact Match has been added.

7. Click **Details**.
8. Select **Duplicate Collections** from the show list and then click **Refresh**.
9. Expand each **CollectionNumber** to view the Suspect and Duplicate records for each duplicate collection.



10. Compare the collections in the Detail view to the output file created.

Dataflow Templates for Matching

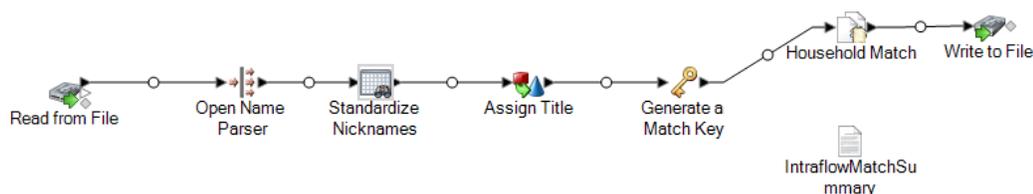
Identifying Members of a Household

This dataflow template demonstrates how to identify members of the same household by comparing information within a single input file and creating an output file of household collections.

Business Scenario

As data steward for a credit card company and you want to analyze your customer database and find out which addresses occur multiple times and under what names so that you can minimize that number of duplicate mailings and credit card offers sent to the same address.

The following dataflow provides a solution to the business scenario:



This dataflow template is available in Enterprise Designer. Go to **File > New > Dataflow > From template** and select **HouseholdRelationships**. This dataflow requires the following products: Advanced Matching, Data Normalization, and Universal Name.

For each record in the input file, this dataflow will do the following:

Read from File

This stage identifies the file name, location, and layout of the file that contains the names you want to parse. The file contains both male and female names.

Open Name Parser

Open Name Parser examines name fields and compares them to name data stored in the Spectrum Technology Platform name database files. Based on the comparison, it parses the name data into First, Middle, and Last name fields, assigns an entity type, and a gender to each name. It also uses pattern recognition in addition to the name data.

Standardize Nicknames

In this template, the Table Lookup stage is named Standardize Nicknames. Standardize Nickname stage looks up first names in the Nicknames.xml database and replaces any nicknames with the more regular form of the nickname. For example, the name Tommy is replaced with Thomas.

Transformer

In this template, the Transformer stage is named Assign Titles. Assign Titles stage uses a custom script to search each row in the data stream output by the Parse Personal Name stage and assign a **TitleOfRespect** value based on the **GenderCode** value.

The custom script is:

```
if (row.get('TitleOfRespect') == '')
{
  if (row.get('GenderCode') == 'M')
    row.set('TitleOfRespect', 'Mr')
  if (row.get('GenderCode') == 'F')
    row.set('TitleOfRespect', 'Ms')
```

Every time the Assign Titles stage encounters **M** in the **GenderCode** field it sets the value for **TitleOfRespect** as **Mr**. Every time the Assign Titles stages encounters **F** in the **GenderCode** field it sets the value of **TitleOfRespect** as **Ms**.

Match Key Generator

The Match Key Generator processes user-defined rules that consist of algorithms and input source fields to generate the match key field. A match key is a non-unique key shared by like records that

identify records as potential duplicates. The match key is used to facilitate the matching process by only comparing records that contain the same match key. A match key is comprised of input fields. Each input field specified has a selected algorithm that is performed on it. The result of each field is then concatenated to create a single match key field.

In this template, two match key fields are defined: SubString (LastName (1:3)) and SubString (PostalCode (1:5)).

For example, if the incoming address was:

FirstName - Fred

LastName - Mertz

PostalCode - 21114-1687

And the rules specified that:

Input Field	Start Position	Length
LastName	1	3
PostalCode	1	5

Then the key, based on the rules and the input data shown above, would be:

Mer21114

Household Match

In this dataflow template the Intraflow Match stage is named Household Match. This stage locates matches between similar data records within a single input stream. Matched records can also be qualified by using non-name/non-address information. The matching engine allows you to create hierarchical rules based on any fields that have been defined or created in other stages.

A stream of records to be matched as well as settings that specify what fields should be compared, how scores should be computed, and generally what constitutes a successful match.

In this template, you create a custom matching rule that compares LastName and AddressLine1. Select the **Generate data for analysis** check box to generate data for the Interflow Summary Report.

Here are some guidelines to follow when creating your matching hierarchy:

- A parent node must be given a unique name. It can not be a field.
- The child field must be a Spectrum Technology Platform data type field, that is, one available through one or more components.

- All children under a parent must use the same logical operators. To combine connectors you must first create intermediate parent nodes.
- Thresholds at the parent node could be higher than the threshold of the children.
- Parent nodes do not have to have a threshold.

Write to File

The template contains one Write to File stage that creates a text file that shows the addresses as a collection of households.

Intraflow Summary Report

The template contains the Intraflow Match Summary Report. After you run the job, expand **Reports** in the Execution Details window, and then click **IntraflowMatchSummary**.

The Intraflow Match Summary Report lists the statistics for the records processed and shows a bar chart that graphically illustrates the record count and overall matching score.

Determining if a Prospect is a Customer

This dataflow template demonstrates how to evaluate prospect data in an input file to customer data in a customer database to determine if a prospect is a customer. This is a service dataflow, meaning that the dataflow can be accessed via the API or web services.

Business Scenario

As a sales executive for an online sales company you want to determine if an online prospect is an existing customer or a new customer.

The following dataflow service provides a solution to the business scenario:



This dataflow template is available in Enterprise Designer. Go to **File > New > Dataflow > From template** and select **ProspectMatching**. This dataflow requires these products: Advanced Matching and Universal Name.

For each record in the input file, this dataflow does the following:

Input

The selected input fields for this template are AddressLine1, City, Name, PostalCode, and StateProvince. AddressLine1 and Name are the fields that are key to the dataflow processing in this template.

Name Parser

In this template, the Name Parser stage is named Parse Personal Name. Parse Personal Name stage examines name fields and compares them to name data stored in the Spectrum Technology Platform name database files. Based on the comparison, it parses the name data into First, Middle, and Last name fields, assigns an entity type, and a gender to each name. It also uses pattern recognition in addition to the name data.

In this template the Parse Personal Name stage is configured as follows.

- Parse personal names is selected and Parse business names is cleared. When you select these options, first names are evaluated for gender, order, and punctuation and no evaluation of business names is performed.
- Gender Determination Source is set to default. For most cases, Default is the best setting for gender determination because it covers a wide variety of names. However, if you are processing names from a specific culture, select that culture. Selecting a specific culture helps ensure that the proper gender is assigned to the names. For example, if you leave Default selected, then the name Jean will be identified as a female name. However, if you select French, it will be identified as a male name.
- Order is set to natural. The name fields are ordered by Title, First Name, Middle Name, Last Name, and Suffix.
- Retain periods is cleared. Any punctuation in the name data is not retained.

Candidate Finder

The Candidate Finder stage is used in combination with the Transactional Match stage.

The Candidate Finder stage obtains the candidate records that will form the set of potential matches that the Transactional Match stage will evaluate. In addition, depending on the format of your data, Candidate Finder may need to parse the name or address of the suspect record, the candidate records, or both.

As part of configuring Candidate Finder, you select the database connection through which the specified query will be executed. You can select any connection configured in Management Console. To connect to a database not listed, configure a connection to that database in Management Console, then close and reopen Candidate Finder to refresh the connection list.

To define the SQL query you can type any valid SQL select statement into the text box on the Candidate Finder Options view. For example, assume you have a table in your database called Customer_Table that has the following columns:

Customer_Table

Cust_Name

Cust_Address

Cust_City

Cust_State

Cust_Zip

Note: You can type any valid SQL select, however, `Select *` is not valid in this control.

To retrieve all the rows from the database, you might construct a query similar to the following:

```
select Cust_Name, Cust_Address, Cust_City, Cust_State, Cust_Zip from
Customer_Table;
```

However, it is unlikely that you would want to match your transaction against all the rows in the database. To return only relevant candidate records, you will want to add a WHERE clause using variable substitution. Variable substitution refers to a special notation that you will use to cause the Candidate Selection engine to replace the variable with the actual data from your suspect record.

To use variable substitution, enclose the field name in braces preceded by a dollar sign using the form **`\${FieldName}**. For example, the following query will return only those records that have a value in `Cust_Zip` that matches the value in `PostalCode` on the suspect record.

```
select Cust_Name, Cust_Address, Cust_City, Cust_State, Cust_Zip
from Customer_Table
where Cust_Zip = `${PostalCode};
```

Next you need to map database columns to stage fields if the column names in your database do not match the Component Field names exactly. If they do match they will be automatically mapped to the corresponding Stage Fields. You will need to use the Selected Fields (columns from the database) to map to the Stage Fields (field names defined in the dataflow).

Again consider the `Customer_Table` from the above example:

Customer_Table

Cust_Name

Cust_Address

Cust_City

Cust_State

Cust_Zip

When you retrieve these records from the database, you need to map the column names to the field names that will be used by the Transactional Match stage and other stages in your dataflow. For example, Cust_Address might be mapped to AddressLine1, and Cust_Zip would be mapped to PostalCode.

1. Select the drop-down list under **Selected Fields** in the candidate Finder Options view. Then, select the database column Cust_Zip.
2. Select the drop-down list under **Stage Fields**. Then, select the field to which you want to map.

For example, if you want to map Cust_Zip to Postal Code, first select Cust_Zip under Selected fields and then select PostalCode on the corresponding Stage Field row.

In addition to mapping fields as described above, you can use special notation in your SQL query to perform the mapping. To do this, you will enter the name of the Stage Field, enclosed in braces, after the column name in your query. When you do this, the selected fields will be automatically mapped to the corresponding stage fields.

An example of this using the query from the previous example follows:

```
select Cust_Name {Name}, Cust_Address {AddressLine1},
       Cust_City {City}, Cust_State {StateProvince},
       Cust_Zip {PostalCode}
from Customer
where Cust_Zip = ${PostalCode};
```

Transactional Match

The Transactional Match stage is used in combination with the Candidate Finder stage.

The Transactional Match stage allows you to match suspect records against potential candidate records that are returned from the Candidate Finder Stage.

Transactional Match uses matching rules to compare the suspect record to all candidate records with the same candidate group number (assigned in Candidate Finder) to identify duplicates. If the candidate record is a duplicate, it is assigned a collection number, the match record type is labeled a Duplicate, and the record is then written out. Any unmatched candidates in the group are assigned a collection number of 0, labeled as Unique and then written out as well.

In this template, you create a custom matching rule that compares LastName and AddressLine1.

Here are some guidelines to follow when creating your matching hierarchy:

- A parent node must be given a unique name. It can not be a field.
- The child field must be a Spectrum Technology Platform data type field, that is, one available through one or more stages.
- All children under a parent must use the same logical operators. To combine connectors you must first create intermediate parent nodes.
- Thresholds at the parent node could be higher than the threshold of the children.
- Parent nodes do not have to have a threshold.

Output

As a service, this template sends all available fields to the output. You can limit the output based on your needs.

5 - Deduplication

In this section

Filtering Out Duplicate Records.....	135
Creating a Best of Breed Record.....	139



Filtering Out Duplicate Records

The simplest way to remove duplicate records is to add a Filter stage to your dataflow after a matching stage. The Filter stage removes records from collections of duplicate records based on the settings you specify.

1. In Enterprise Designer, create a dataflow that identifies duplicate records through matching.

Matching is the first step in deduplication because you need to identify records that are similar, such as records that have the same account number or name. See the following topics for instructions on creating a dataflow that matches records.

[Matching Records from a Single Source](#) on page 80

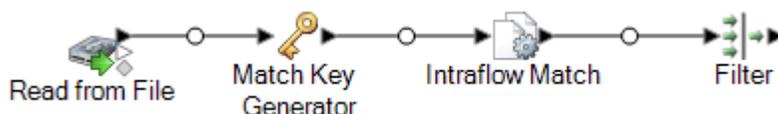
[Matching Records from One Source to Another Source](#) on page 86

[Matching Records Against a Database](#) on page 98

Note: You only need to build the dataflow to the point where it reads data and performs matching with an Interflow Match, Intraflow Match, or Transactional Match stage. Once you have created a dataflow to this point, continue with the following steps.

2. Once you have defined a dataflow that reads data and matches records, drag a Filter stage to the canvas and connect it to the stage that performs the matching (Interflow Match, Intraflow Match, or Transactional Match).

For example, if your dataflow reads data from a file and performs matching with Intraflow Match, your dataflow would look like this after adding a Filter stage:



3. Double-click the Filter stage on the canvas.
4. In the **Group by** field, select **CollectionNumber**.
5. Leave the option **Limit number of returned duplicate records** selected and the value set to 1. These are the default settings.
6. Decide if you want to keep the first record in each collection, or if you want to define a rule to choose which record from each collection to keep. If you want to keep the first record in each collection, skip this step. If you want to define a rule, in the rule tree, select **Rules** then follow these steps:
 - a) Click **Add Rule**.

Records in each group are evaluated to see if they meet the rules you define here. If a record meets the rule, it is the surviving record and the other records in the group are discarded.

b) Define a rule to identify the record from each group to retain.

Use the following options to define a rule:

Option	Description
Field name	Specifies the name of the dataflow field whose value you want to evaluate to determine whether to filter the record.
Field Type	Specifies the type of data in the field. One of the following: Non-Numeric Choose this option if the field contains non-numeric data (for example, string data). Numeric Choose this option if the field contains numeric data (for example, double, float, and so on).

Option	Description
Operator	Specifies the type of comparison you want to use to evaluate the field. One of the following: <ul style="list-style-type: none"> Contains Determines if the field contains the value specified. For example, "sailboat" contains the value "boat". Equal Determines if the field contains the exact value specified. Greater Than Determines if the field value is greater than the value specified. This operation only works on numeric fields. Greater Than Or Equal To Determines if the field value is greater than or equal to the value specified. This operation only works on numeric fields. Highest Compares the field's value for all the records group and determines which record has the highest value in the field. For example, if the fields in the group contain values of 10, 20, 30, and 100, the record with the field value 100 would be selected. This operation only works on numeric fields. If multiple records are tied for the longest value, one record is selected. Is Empty Determines if the field contains no value. Is Not Empty Determines if the field contains any value. Less Than Determines if the field value is less than the value specified. This operation only works on numeric fields. Less Than Or Equal To Determines if the field value is less than or equal to the value specified. This operation only works on numeric fields. Longest Compares the field's value for all the records group and determines which record has the longest (in bytes) value in the field. For example, if the group contains the values "Mike" and "Michael", the record with the value "Michael" would be selected. If multiple records are tied for the longest value, one record is selected. Lowest Compares the field's value for all the records group and determines which record has the lowest value in the field. For example, if the fields in the group contain values of 10, 20, 30, and 100, the record with the field value 10 would be selected. This operation only works on numeric fields. If multiple records are tied for the longest value, one record is selected. Most Common Determines if the field value contains the value that occurs most frequently in this field among the records in the group. If two or more values are most common, no action is taken. Not Equal Determines if the field value is not the same as the value specified.

Option	Description
Value type	<p>Specifies the type of value you want to compare to the field's value. One of the following:</p> <p>Note: This option is not available if you select the operator Highest, Lowest, or Longest.</p> <p>Field Choose this option if you want to compare another dataflow field's value to the field.</p> <p>String Choose this option if you want to compare the field to a specific value.</p>
Value	<p>Specifies the value to compare to the field's value. If you selected Field in the Field type field, select a dataflow field. If you selected String in the Value type field, type the value you want to use in the comparison.</p> <p>Note: This option is not available if you select the operator Highest, Lowest, or Longest.</p>

c) Click **OK**.

You have now configured Filter with one rule. You can add additional rules if needed.

7. Click **OK** to close the **Filter Options** window.
8. Drag a sink stage onto the canvas and connect it to the Filter stage.

For example, if you were using a Write to File sink stage your dataflow would look like this:



9. Double-click the sink stage and configure it.

For information on configuring sink stages, see the *Dataflow Designer's Guide*.

You now have a dataflow that identifies matching records and removes all but one record for each group of duplicates, resulting in an output file that contains deduplicated data.

Creating a Best of Breed Record

To eliminate duplicate records from your data, you may choose to merge data from groups of duplicate records into a single "best of breed" record. This approach is useful when each duplicate record contains data of the same type (for example, phone numbers or names) and you want to preserve the best data from each record in the surviving record.

This procedure describes how create a dataflow that merges duplicate records into a best of breed record.

1. In Enterprise Designer, create a dataflow that identifies duplicate records through matching.

Matching is the first step in deduplication because you need to identify records that are similar, such as records that have the same account number or name. See the following topics for instructions on creating a dataflow that matches records.

[Matching Records from a Single Source](#) on page 80

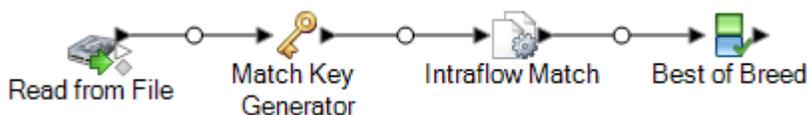
[Matching Records from One Source to Another Source](#) on page 86

[Matching Records Against a Database](#) on page 98

Note: You only need to build the dataflow to the point where it reads data and performs matching with an Interflow Match, Intraflow Match, or Transactional Match stage. Once you have created a dataflow to this point, continue with the following steps.

2. Once you have defined a dataflow that reads data and matches records, drag a Best of Breed stage to the canvas and connect it to the stage that performs the matching (Interflow Match, Intraflow Match, or Transactional Match).

For example, if your dataflow reads data from a file and performs matching with Intraflow Match, your dataflow would look like this after adding a Best of Breed stage:



3. Double-click the Best of Breed stage on the canvas.
4. In the **Group by** field, select **CollectionNumber**.
5. Under **Best of Breed Settings**, select **Rules** in the conditions tree.
6. Click **Add Rule**.

Records in each group are evaluated to see if they meet the rules you define here. If a record matches a rule, its data may be copied to the best of breed record, depending on how you configure the actions associated with the rule. You will define actions later.

7. Define a rule that a duplicate record must meet in order for its data to be copied to the best of breed record.

Use the following options to define a rule:

Option	Description
Field name	Specifies the name of the dataflow field whose value you want to evaluate to determine if the condition is met and the associated actions should be taken.
Field Type	Specifies the type of data in the field. One of the following: Non-Numeric Choose this option if the field contains non-numeric data (for example, string data). Numeric Choose this option if the field contains numeric data (for example, double, float, and so on).

Option	Description
Operator	Specifies the type of comparison you want to use to evaluate the field. One of the following:
Contains	Determines if the field contains the value specified. For example, "sailboat" contains the value "boat".
Equal	Determines if the field contains the exact value specified.
Greater Than	Determines if the field value is greater than the value specified. This operation only works on numeric fields.
Greater Than Or Equal To	Determines if the field value is greater than or equal to the value specified. This operation only works on numeric fields.
Highest	Compares the field's value for all the records group and determines which record has the highest value in the field. For example, if the fields in the group contain values of 10, 20, 30, and 100, the record with the field value 100 would be selected. This operation only works on numeric fields. If multiple records are tied for the longest value, one record is selected.
Is Empty	Determines if the field contains no value.
Is Not Empty	Determines if the field contains any value.
Less Than	Determines if the field value is less than the value specified. This operation only works on numeric fields.
Less Than Or Equal To	Determines if the field value is less than or equal to the value specified. This operation only works on numeric fields.
Longest	Compares the field's value for all the records group and determines which record has the longest (in bytes) value in the field. For example, if the group contains the values "Mike" and "Michael", the record with the value "Michael" would be selected. If multiple records are tied for the longest value, one record is selected.
Lowest	Compares the field's value for all the records group and determines which record has the lowest value in the field. For example, if the fields in the group contain values of 10, 20, 30, and 100, the record with the field value 10 would be selected. This operation only works on numeric fields. If multiple records are tied for the longest value, one record is selected.
Most Common	Determines if the field value contains the value that occurs most frequently in this field among the records in the group. If two or more values are most common, no action is taken.
Not Equal	Determines if the field value is not the same as the value specified.

Option	Description
Value type	<p>Specifies the type of value you want to compare to the field's value. One of the following:</p> <p>Note: This option is not available if you select the operator Highest, Lowest, or Longest.</p> <p>Field Choose this option if you want to compare another dataflow field's value to the field.</p> <p>String Choose this option if you want to compare the field to a specific value.</p>
Value	<p>Specifies the value to compare to the field's value. If you selected Field in the Field type field, select a dataflow field. If you selected String in the Value type field, type the value you want to use in the comparison.</p> <p>Note: This option is not available if you select the operator Highest, Lowest, or Longest.</p>

8. Click **OK**.
9. Click the **Actions** node in the tree.
10. Click **Add Action**.
11. Specify the data to copy to the best of breed record if the record meets the criteria you defined in the rule.

Option	Description
Source type	<p>Specifies the type of data to copy to the best of breed record. One of the following.</p> <p>Field Choose this option if you want to copy a value from a field to the best of breed record.</p> <p>String Choose this option if you want to copy a constant value to the best of breed record.</p>
Source data	<p>Specifies the data to copy to the best of breed record. If the source type is Field, select the field whose value you want to copy to the destination field. If the source type is String, specify a constant value to copy to the destination field.</p>
Destination	<p>Specifies the field in the best of breed record to which you want to copy the data specified in the Source data field.</p>

Option	Description
Accumulate source data	<p>If the data in the Source data field is numeric data, you can enable this option to combine the source data for all duplicate records and put the total value in the best of breed record.</p> <p>For example, if there were three duplicate records in the group and they contained these values in the Deposits field:</p> <p>100.00 20.00 5.00</p> <p>Then all three values would be combined and the total value, 125.00, would be put in the best of breed record's Deposits field.</p>

- Click **OK**.

You have now configured Best of Breed with one rule and one action. You can add additional rules and actions if needed.

- Click **OK** to close the **Best of Breed Options** window.
- Drag a sink stage onto the canvas and connect it to the Best of Breed stage.

For example, if you were using a Write to File sink stage your dataflow would look like this:



- Double-click the sink stage and configure it.

For information on configuring sink stages, see the *Dataflow Designer's Guide*.

You now have a dataflow that identifies matching records and merges records within a collection into a single best of breed record.

6 - Exception Records

In this section

Designing a Dataflow to Handle Exceptions.....	145
Designing a Dataflow for Real-Time Revalidation.....	146



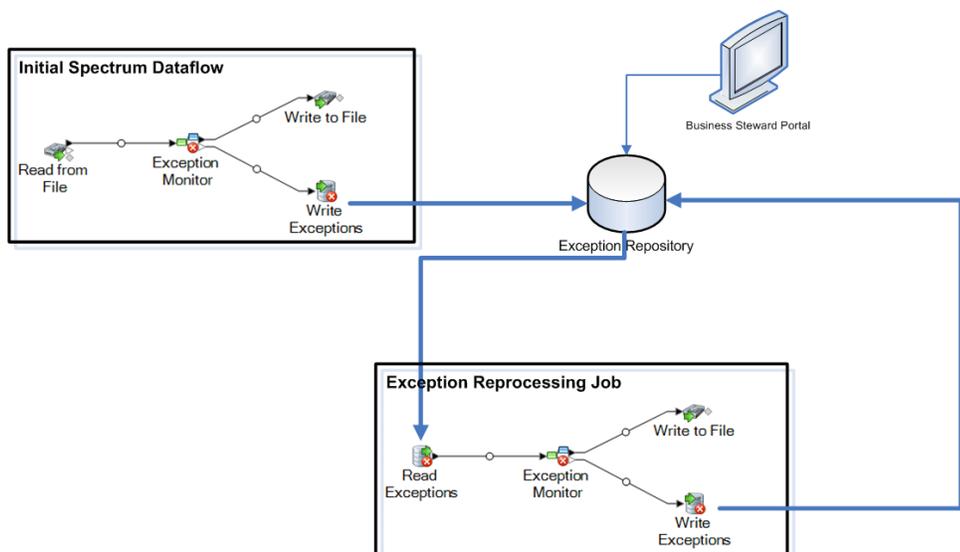
Designing a Dataflow to Handle Exceptions

If you have licensed Data Stewardship, you can include an exception management process in your dataflows. The basic building blocks of an exception management process are:

- An initial dataflow that performs a data quality process, such as record deduplication, address validation, or geocoding.
- An Exception Monitor stage that identifies records that could not be processed.
- A Write Exceptions stage that takes the exception records identified by the Exception Monitor stage and writes them to the exception repository for manual review.
- The Data Stewardship Portal, a browser-based tool, which allows you to review and edit exception records. After records are edited, they are resolved and reprocessed.
- An exception reprocessing job that uses the Read Exceptions stage to read approved records from the exception repository into the job. The job then attempts to reprocess the corrected records, typically using the same logic as the original dataflow. The Exception Monitor stage once again checks for exceptions. The Write Exceptions stage then sends exceptions back to the exception repository for additional review.

Note: Do not place other stages between Exception Monitor and Write Exceptions stages in a dataflow; doing so could impact the configuration of exceptions in the Data Stewardship Portal.

Here is an example scenario that helps illustrate a basic exception management implementation:



In this example, there are two dataflows: the initial dataflow, which evaluates the input records' postal code data, and the exception reprocessing job, which takes the edited exceptions and verifies that the records now contain valid postal code data.

In both dataflows there is an Exception Monitor stage. This stage contains the conditions you want to use to determine if a record should be routed for manual review. These conditions consist of one or more expressions, such as `PostalCode is empty`, which means any record not containing a postal code would be considered an exception and would be routed to the Write Exceptions stage and written to the exception repository. For more information, see [Exception Monitor](#) on page 310.

Any records that the Exception Monitor identifies as exceptions are routed to an exception repository using the Write Exceptions stage. Data stewards review the exceptions in the repository using the Data Stewardship Portal, a browser-based tool for viewing and modifying exception records. Using our example, the data steward could use the Data Stewardship Portal **Editor** to manually add postal codes to the exception records and mark them as "Approved".

Once a record is marked resolved in the Data Stewardship Portal, the record is available to be read back into a Spectrum Technology Platform dataflow. This is accomplished by using a Read Exceptions stage. Should any records still result in an exception, they are again written to the exception repository for review by a data steward.

To determine the best approach for your situation, consider these questions:

- **How do you want to identify exception records?** The Exception Monitor stage can evaluate any field's value or any combination of fields to determine if a record is an exception. You should analyze the results you are currently getting with your dataflow to determine how you want to identify exceptions. You may want to identify records in the middle range of the data quality continuum, and not those that were clearly validated or clearly failed.
- **Do you want edited and approved exception records re-processed using the same logic as was used in the original dataflow?** If so you may want to use a subflow to create reusable business logic. For example, the subflow could be used in an initial dataflow that performs address validation and in an exception reprocessing job that re-processes the corrected records to verify the corrections. You can then use different source and sink stages between the two. The initial dataflow might contain a Read from DB stage that takes data from your customer database for processing. The exception reprocessing job would contain a Read Exceptions stage that takes the edited and approved exception records from the exception repository.
- **Do you want to reprocess corrected and approved exceptions on a predefined schedule?** If so you can schedule your reprocessing job using Scheduling in the Management Console.

Designing a Dataflow for Real-Time Revalidation

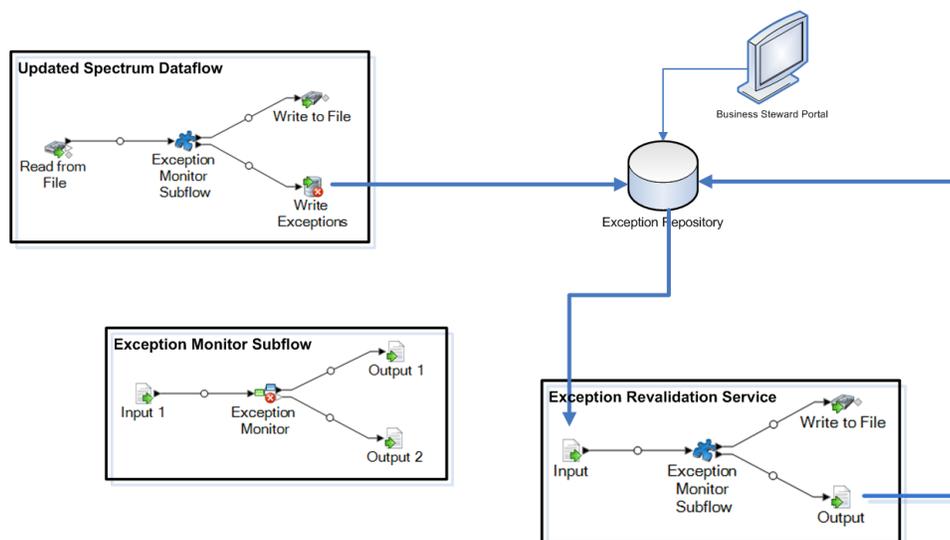
If you are using exception management in your dataflow, you can use the revalidation feature to rerun exception records through the validation process after they have been corrected in the Data Stewardship Portal. This enables you to determine if the change you made causes the record to

process successfully in a real-time manner. You do not need to wait until the Read Exceptions batch job runs again to see the result.

The basic building blocks of a revalidation environment are:

- A job or a service that reuses or contains an exposed subflow. It must also contain an input source, the subflow stage that processes the input, a Write Exceptions stage, and an output sink for successfully processed records.
- An exposed subflow containing an Exception Monitor stage that points to a revalidation service and is configured for revalidation, including designating whether revalidated records should be reprocessed or approved.
- An exposed service that also reuses or contains the exposed subflow. It processes records that were edited, saved, and sent for revalidation in the Data Stewardship Portal.

Here is an example scenario that helps illustrate a revalidation implementation:



In this example, there are three dataflows: a job, a subflow, and a service. The job runs input data through the subflow. The subflow contains an Exception Monitor stage, which determines if a record should be routed for manual review. That means any records with no data in the PostalCode field would be considered an exception and would be routed to the Write Exceptions stage; these exceptions are what appears in the Data Stewardship Portal. Records with anything else in that field would be routed to the Write to File stage.

Note: If your dataflow is also being configured to use best of breed functionality, you will need to manually add and expose the CollectionRecordType field in the revalidation Exception Monitor stage/subflow and the service itself. See [Write Exceptions](#) options and [Creating a Best of Breed Record](#) on page 190 for more information on best of breed functionality.

The exception revalidation service that you designated when configuring the Exception Monitor stage is called when you edit one or more records from the exception repository on the Data Stewardship Portal **Editor** page and click **Save**. Like the job, the service contains the exception monitor subflow that uses the same business logic to reprocess the records. If the records fail one or more conditions

set in the Exception Monitor stage, the exceptions will be updated in the repository. If the records pass the conditions set in the Exception Monitor stage, one of two actions will occur, depending on the selection made for **Action after revalidation**:

Reprocess records	Records will be deleted from the repository and reprocessed.
Resolve records	The approved records are retained in the repository and their status changed to Resolved.

To step through the procedure, see [Create a real-time validation scenario](#) on page 148.

Create a real-time validation scenario

Follow these steps to create and use a real-time revalidation scenario:

1. Open or create a job or service dataflow that contains an Exception Monitor stage, an input source (such as a Read from File or Input stage), an output sink (such as a Write to File or Output stage), and a Write Exceptions stage.
2. Convert the Exception Monitor stage to a subflow and map the input and output fields to match those in the initial dataflow. Be sure to include the ExceptionMetadata field for the input source as well as the output stage that populates the Write Exceptions stage in the job. Expose the subflow so it can be used by the job and service.
3. Create a service that contains an Input stage, the subflow you created in step 2 on page 148, an Output stage, and an output sink (such as a Write to File or Write to DB stage). Map the input and output fields to match those in the initial dataflow; be sure to include the ExceptionMetadata field for the Input stage as well as the Output stage. Expose the service so it can be used by the subflow.
4. Return to the subflow and open the Configuration tab of the Exception Monitor stage. Select the revalidation service you created in step 3 on page 148 and specify which action to take after revalidation. Save and expose the subflow again.
5. Return to the service, where a message will appear, notifying you of changes to the subflow and saying that the service will be refreshed. Click **OK**, then save and expose the service again.
6. Return to the initial job or service, where a message will appear, notifying you of changes to the subflow and saying that the dataflow will be refreshed. Click **OK**, then save the dataflow.
7. Run the job.

Note: Even if you have run the initial job or service before, you must run it again after creating the revalidation scenario to populate the repository with records that are eligible for revalidation. You can identify whether records in the Data Stewardship Portal **Editor** are eligible for revalidation because the **Save** button will be active for those records.

7 - Lookup Tables

In this section

Introduction to Lookup Tables.....	150
Data Normalization Tables.....	150
Universal Name Tables.....	155
Viewing the Contents of a Lookup Table.....	157
Adding a Term to a Lookup Table.....	158
Removing a Term from a Lookup Table.....	158
Modifying the Standardized Form of a Term.....	159
Reverting Table Customizations.....	159
Creating a Lookup Table.....	160
Importing Data.....	160



Introduction to Lookup Tables

A lookup table is a table of key/value pairs used by Spectrum Technology Platform stages to standardize data by performing token replacement. To modify the contents of the lookup tables used in Advanced Transformer, Open Parser, and Table Lookup, use the Table Management tool in Enterprise Designer.

Data Normalization Tables

Advanced Transformer Tables

Advanced Transformer uses the following tables to identify terms. Use Table Management to create new tables or to modify existing ones. For more information, see [Introduction to Lookup Tables](#) on page 150.

- Aeronautical Abbreviations
- All Acronyms Initialism
- Business Names Abbreviations
- Canadian Territory Abbreviations
- Computing/IT Abbreviations
- Delimiters
- German Companies
- Fortune 1000
- Geographic Directional Abbreviations
- Global Sentry Noise Terms
- Global Sentry Sanctioned Countries
- Government Agencies Abbreviations
- IATA Airline Designator
- IATA Airline Designator Country
- Legal Abbreviations
- Medical Abbreviations
- Medical Organizations Acronyms
- Military Abbreviations
- Nicknames

- Secondary Unit Abbreviations
- Secondary Unit Reverse
- Singapore Abbreviations
- Spanish Abbreviations
- Spanish Directional Abbreviations
- Spanish Street Suffix Abbreviations
- State Name Abbreviations
- State Name Reverse
- Street Suffix Abbreviations
- Street Suffix Reverse
- Subsidiary to Parent
- U.S. Army Acronyms
- U.S. Navy Acronyms

Open Parser Tables

Open Parser uses the following tables to identify terms. Use Table Management to create new tables or to modify existing ones. For more information, see [Introduction to Lookup Tables](#) on page 150.

Base Tables

Base tables are provided with the Data Normalization installation package.

- Account Descriptions
- Companies
- Company Conjunctions
- Company Prepositions
- Company Suffixes
- Company Terms
- Conjunctions
- Family Name Prefixes
- Family Names
- General Suffixes
- German Companies
- Given Names
- Maturity Suffixes
- Spanish Given Names
- Spanish Family Names
- Titles

Core Name Tables

Core Names tables are not provided with the Data Normalization installation package and thus require an additional license. For more information, contact your account executive.

Core Names tables must be loaded using the Data Normalization database load utility. For instructions, see the *Spectrum Technology Platform Installation Guide*.

- Enhanced Family Names
- Enhanced Given Names

Company Name Tables

Company Names tables are not provided with the Data Normalization installation package and thus require an additional license. For more information, contact your account executive.

Company Names tables must be loaded using the Data Normalization database load utility. For instructions, see the *Spectrum Technology Platform Installation Guide*.

- Companies - Americas
- Companies - Asia Pacific
- Companies - EMEA
- Company Articles
- Company Conjunctions

Arabic Plus Pack Tables

Arabic Plus Pack tables are not provided with the Data Normalization installation package and thus require an additional license. For more information, contact your account executive.

Arabic Plus Pack tables must be loaded using the Data Normalization database load utility. For instructions, see the *Spectrum Technology Platform Installation Guide*.

- Arabic Family Names (Arabic)
- Arabic Family Names (Romanized)
- Arabic Given Names (Arabic)
- Arabic Given Names (Romanized)

Asian Plus Pack Tables

Asian Plus Pack tables are not provided with the Data Normalization installation package require an additional license. For more information, contact your account executive.

Asian Plus Pack tables must be loaded using the Data Normalization database load utility. For instructions, see the *Spectrum Technology Platform Installation Guide*.

- Chinese Family Names (Native)
- Chinese Family Names (Romanized)
- Chinese Given Names (Native)
- Chinese Given Names (Romanized)
- Korean Family Names (Native)
- Korean Family Names (Romanized)
- Korean Given Names (Native)
- Korean Given Names (Romanized)
- Japanese Family Names (Kana)
- Japanese Family Names (Kanji)
- Japanese Family Names (Romanized)
- Japanese Given Names (Kana)
- Japanese Given Names (Kanji)
- Japanese Given Names (Romanized)

Table Lookup Tables

Table Lookup uses the following tables to identify terms. Use Table Management to create new tables or to modify existing ones. For more information, see [Introduction to Lookup Tables](#) on page 150.

Base Tables

Base tables are provided with the Data Normalization installation package.

- Aeronautical Abbreviations
- All Acronyms Initialism
- Business Names Abbreviations
- Canadian Territory Abbreviations
- Computing/IT Abbreviations
- EU Acronyms
- Fortune 1000
- French Abbreviations
- French Arrondissement to Department Number
- French Commune to Postal Code
- French Department to Region
- French Department Number to Department
- Gender Codes
- Geographic Directional Abbreviations
- German Acronyms

- German City to State Code
- German Area Code to City
- German District to State Code
- German State Abbreviations
- Global Sentry Sanctioned Countries
- Government Agencies Abbreviations
- IATA Airline Designator
- IATA Airline Designator Country
- Legal Abbreviations
- Medical Abbreviations
- Medical Organizations Acronyms
- Military Abbreviations
- Nicknames
- Secondary Unit Abbreviations
- Secondary Unit Reverse
- Singapore Abbreviations
- Spanish Abbreviations
- Spanish Directional Abbreviations
- Spanish Street Suffix Abbreviations
- State Name Abbreviations
- State Name Reverse
- Street Suffix Abbreviations
- Street Suffix Reverse
- Subsidiary to Parent
- U.K. Town to Postcode Area
- U.K. Dialing Code Prefixes
- U.K. Dialing Codes to Town
- U.K. Postcode Area to Town
- U.S. Army Acronyms
- U.S. Navy Acronyms
- ZREPLACE (Used by the SAP Module for French address validation)

Core Names

Core Names tables require an additional license. For more information, contact your account executive.

Core Names tables must be loaded using the Data Normalization database load utility. For instructions, see the *Spectrum Technology Platform Installation Guide*.

- Enhanced Family Names Ethnicity
- Enhanced Gender Codes
- Enhanced Given Names Ethnicity

Arabic Plus Pack

Arabic Plus Pack tables require an additional license. For more information, contact your account executive.

Arabic Plus Pack tables must be loaded using the Data Normalization database load utility. For instructions, see the *Spectrum Technology Platform Installation Guide*.

- Arabic Family Names Ethnicity (Arabic)
- Arabic Family Names Ethnicity (Romanized)
- Arabic Gender Codes (Arabic)
- Arabic Gender Codes (Romanized)
- Arabic Given Names Ethnicity (Arabic)
- Arabic Given Names Ethnicity (Romanized)

Asian Plus Pack

Asian Plus Pack tables require an additional license. For more information, contact your account executive.

Asian Plus Pack tables must be loaded using the Data Normalization database load utility. For instructions, see the *Spectrum Technology Platform Installation Guide*.

- CJK Family Names Ethnicity (Native)
- CJK Family Names Ethnicity (Romanized)
- CJK Given Names Ethnicity (Native)
- CJK Given Names Ethnicity (Romanized)
- Japanese Gender Codes (Kana)
- Japanese Gender Codes (Kanji)
- Japanese Gender Codes (Romanized)

Universal Name Tables

Name Variant Finder Tables

The Name Variant Finder stage uses the following tables. Each table requires a separate license.

- Arabic Plus Pack: `g1-cdq-cjki-arabic-<date>.jar`
- Asian Plus Pack - Chinese: `g1-cdq-cjki-chinese-<date>.jar`
- Asian Plus Pack - Japanese: `g1-cdq-cjki-japanese-<date>.jar`

- Asian Plus Pack - Korean: g1-cdq-cjki-korean-<date>.jar
- Core Names Database: g1-cdq-nomino-base-<date>.jar

Open Name Parser Tables

Open Name Parser uses the following tables to identify terms. Use Table Management to create new tables or to modify existing ones. For more information, see [Introduction to Lookup Tables](#) on page 150.

Base Tables

Base tables are provided with the Universal Name installation package.

- Account Descriptions
- Company Conjunctions
- Conjunctions
- Family Name Prefixes
- Family Names
- General Suffixes
- Given Names
- Maturity Suffixes
- Spanish Given Names
- Spanish Family Names
- Titles

Core Name Tables

Core name tables are not provided with the Universal Name installation package and thus require an additional license.

- Enhanced Family Names
- Enhanced Given Names

Company Name Tables

The following company name tables are provided with the Universal Name installation package.

- Account Descriptions
- Companies
- Company Articles
- Company Conjunctions
- Company Prepositions
- Company Suffixes
- Company Terms

- Conjunctions

The following company name tables are not provided with the Universal Name installation package and thus require an additional license.

- Companies - Americas
- Companies - Asia Pacific
- Companies - EMEA

Asian Plus Pack Tables

Asian Plus Pack tables are not provided with the Universal Name installation package and thus require an additional license.

- Japanese Family Names (Kana)
- Japanese Family Names (Kanji)
- Japanese Family Names (Romanized)
- Japanese Given Names (Kana)
- Japanese Given Names (Kanji)
- Japanese Given Names (Romanized)
- Japanese Titles

Viewing the Contents of a Lookup Table

You can view the contents of a lookup table by using the Table Management in Enterprise Designer.

1. In Enterprise Designer, select **Tools > Table Management**.
2. In the **Type** field, select the stage whose lookup table you want to view.
3. In the **Name** field, select the table you want to view.
4. You can use the following options to change how the table is displayed:

Option	Description
Find a specific term	In the Starts with field, type the term you want to find then click Refresh .
Page through the table	Click the forward and back icons to the right of the Refresh button.
Change the number of terms displayed per page	Change the value in the Items per page field.

Option	Description
View all the lookup terms for each standardized term in a Table Lookup table	In the View by field select Standardized Term (Grouping) . This option is only available for Table Lookup tables

Adding a Term to a Lookup Table

If you find that your data has terms that are not included in the lookup table and you want to add the term to a lookup table, follow this procedure.

1. In Enterprise Designer, select **Tools > Table Management**.
2. In the **Type** field, select the stage whose lookup table you want to modify.
3. In the **Name** field, select the table to which you want to add a term.
4. Click **Add**.
5. In the **Lookup Term** field, type the term that exists in your data. This is the lookup key that will be used.
6. For Table Lookup tables, in the **Standardized Term** field enter the term you want to be the replacement for the lookup term in your dataflow.

For example, if you want to change the term PS to Precisely Software, you would enter PS as the lookup term, and Precisely Software as the standardized term.

7. For Table Lookup tables, select the **Override existing term** check box if this term already exists in the table and you want to replace it with the value you typed in step 5.
8. Click **Add**.

Removing a Term from a Lookup Table

To remove a term from a lookup table:

1. In Enterprise Designer, select **Tools > Table Management**.
2. Select the term and click **Remove**.
3. Click **Yes** to remove the table term.

Modifying the Standardized Form of a Term

For tables used by Table Lookup to standardize terms, you can change the standardized form of a term. For example, if you have a table where you have the lookup terms PS and PS Inc, and the standardized term is Precisely Software, and you want to change the standardized form to Precisely Software Inc, you could do this by following this procedure.

1. In Enterprise Designer, select **Tools > Table Management**.
2. In the **Type** field, select **Table Lookup**.
3. In the **Name** field select the table you want to modify.
4. Select the term you want to modify and click **Modify**.

Tip: If there are multiple lookup terms for a standardized term, you can easily modify all lookup terms to use the new standardized term by selecting **View by Standardized Term (Grouping)** in the **View by** field, selecting the group, and clicking **Modify**.

5. Type a new value in the **Standardized Term** field.
6. Click **OK**.

Reverting Table Customizations

If you make modifications to a table you can revert the table to its original state. To revert table customizations:

1. In Enterprise Designer, select **Tools > Table Management**.
2. Select the table you want to revert.
3. Click **Revert**.

The **Revert** window displays. It lists all of the added, removed, and modified terms.

4. Select the **Revert** check box for each table entry you want to revert. You can also click **Select All** or **Deselect All** to select or clear all of the **Revert** check boxes.
5. Click **OK**.

Creating a Lookup Table

Advanced Matching, Data Normalization, and Universal Name come with a variety of tables that can be used for a wide range of term replacement or standardization processes. However, if these tables do not meet your needs, you can create your own table of lookup terms to use with Advanced Transformer, Open Parser, or Table Lookup. To create a table, follow this procedure.

1. In Enterprise Designer, select **Tools > Table Management**.
2. In the **Type** field, select the stage for which you want to create a lookup table.
3. Click **New**. The **Add Table** dialog box displays.
4. In the **Table name** field, enter a name for the new table.
5. If you want a new, blank table of the selected type, leave **Copy from** set to **None**. If you want the new table to be populated from an existing table, select a table name from the **Copy from** list.
6. Click **OK**.

For information about adding table items to your new table, see [Adding a Term to a Lookup Table](#) on page 158.

Importing Data

Importing Data Into a Lookup Table

You can import data from a file into a lookup table for use with Advanced Transformer, Open Parser, or Table Lookup. In order to be able to import data from a file into a lookup table, the file must meet these requirements:

- Must be UTF-8 encoded.
- Must be a delimited file. Supported delimiter characters are comma (,), semicolon (;), pipe (|), and tab (\t).
- Fields with embedded delimiters must be start and end with double quotes, for example "1,a","2,b","3,c".
- A literal quote in a field starting and ending with double quote must have two quotes, for example "2"" feet".

To import data from a file into a lookup table:

1. In Enterprise Designer, select **Tools > Table Management**.
2. Select the table into which you want to import the data. Or, create a new table. For instructions on creating a table, see [Creating a Lookup Table](#) on page 160.
3. Click **Import**.
4. Click **Browse** and select the file that contains the data you want to import.
5. Click **Open**. A preview of the data in the imported file displays in Preview File.
6. You can select columns from a user-defined table and map to that in the existing table. For example, assume there are two columns in the user-defined table that you want to import. It has column1 and column2. The column list would show column1 and column2. You could select the column2 to map to a lookup term and select the column1 to map to a standardized term.
7. Select **Import only new terms** to import only new records from the user-defined table or **Overwrite existing terms** to import all records of the selected columns.
8. Click **OK**.

Using Advanced Import

The Advanced Import function allows you to selectively import data into lookup tables used by Advanced Transformer, Table Lookup, and Open Parser. Use Advanced Import to:

- Extract terms from a selected column in a delimited, user-defined file.
- Extract single-word terms (tokens) from a selected column in a delimited user-defined file. When you extract tokens, you can identify the number of times that the terms occurs for a given column in the file and create groupings for related terms and add them to the table.

The file that contains the data you want to import must meet these requirements:

- Must be UTF-8 encoded.
- Must be a delimited file. Supported delimiter characters are comma (,), semicolon (;), pipe (|), and tab (\t).
- Fields with embedded delimiters must be start and end with double quotes, for example "1,a","2,b","3,c".
- A literal quote in a field starting and ending with double quote must have two quotes, for example "2"" feet".

1. In Enterprise Designer, select **Tools > Table Management**.
2. Select the table into which you want to import data.
3. Click **Adv Import**.
4. Click **Browse** and select the file that you want to import.
5. Click **Open**.
6. Select a table column from the **Column** list. The sample data shows the frequency of occurrence for each term listed in the user-defined table. Frequency is only displayed for terms that are not yet in the existing table.

7. To view terms as single words, select **Separate into single-word terms**.
8. For Advanced Transformer and Open Parser tables:
 - a) Select a term from the list on the left.
 - b) Click the right arrow to add the term to the list on the right. Click the left arrow to delete a selected term from the table list.
 - c) Click **OK** to save the changes to the table.
9. For Table Lookup tables:
 - a) Click  to add a table grouping.
 - b) Click **New**.
 - c) Type a new term and then click **Add**. Continue adding terms until finished and then click **Close**.
 - d) Select a term from the list and then click Add. Continue adding terms until finished and then click **Close**. The new terms are added to the terms list on the right.
 - e) Select a term on the left and then click the right arrow to add the term to the selected grouping. Click the left arrow to delete a term from one of the groupings.
 - f) To modify a term, select it from the list on the right and then click .
 - g) To delete a term, select it from the list on the right and then click .
 - h) Click **OK** to save the changes to the table.

8 - Data Stewardship Portal

In this section

Introduction.....	164
Accessing the Data Stewardship Portal.....	164
The Dashboard page.....	166
The Data Stewardship Portal Editor.....	167
The Manage page.....	193
The Data Quality Page.....	198
How to Video - Approval Flow in Data Stewardship.....	199
Introduction.....	199
Accessing the Data Stewardship Portal.....	200
The Dashboard page.....	201
The Data Stewardship Portal Editor.....	202
The Manage page.....	229
The Data Quality Page.....	234
How to Video - Approval Flow in Data Stewardship.....	235
Data Stewardship Settings.....	235



Introduction

The Data Stewardship Portal provides tools for reviewing, modifying, and approving records that failed based on conditions in the Exception Monitor stage. Use the Data Stewardship Portal to manually enter correct or additional data in a record. Typical exceptions are:

- Address verification failures
- Geocoding failures
- Low-confidence matches
- Merge/consolidation decisions

For example, if a customer record fails an address validation process, you could use the search tools to conduct research and determine the customer's address, then modify the record so that it contains the correct address. The modified record could then be approved and reprocessed by Spectrum Technology Platform, sent to another data validation or enrichment process, or written to a database, depending on your configuration. You could also use the Data Stewardship Portal to add information that was not in the original record.

Data Stewardship Portal consists of the following components.

- **Dashboard**—The Dashboard displays the number of records assigned to the user, counts of resolved records that are still assigned to the user, counts of records that are still unresolved, and the daily performance (number of records accepted/rejected). Counts are broken down by Dataflow and Label.
- **Editor**—The editor allows you to modify records and change status of records within review cycles.
- **Manage**—A browser-based tool that enables you to review and manage record activity for all users, including reassigning records from one user to another. Summary and detail charts show the number of records with resolved and unresolved status, and the number of records assigned to each user.
- **Data Quality**—A browser-based tool that provides information on trends for records and enables you to identify key performance indicators and send notifications when certain conditions are met.

Accessing the Data Stewardship Portal

Before you perform this procedure, you must have a user name and password required to access the Spectrum server on your network. Login credentials are normally provided by the Spectrum Technology Platform administrator.

You can open Spectrum and the Data Stewardship Portal using either the Windows Start menu or a web browser.

- **Use the Windows Start menu to access the Data Stewardship Portal** on page 165
- **Access the Data Stewardship Portal from the Spectrum Home page** on page 165
- **Access the Data Stewardship Portal directly** on page 165

All of these methods open the Data Stewardship Portal dashboard.

Use the Windows Start menu to access the Data Stewardship Portal

You can use this Windows Start menu if you are working on the computer that is running the Spectrum Technology Platform server.

1. Click **Start > Precisely > Welcome Page**.
2. Click **Spectrum Data Quality > Data Stewardship Portal > Open the Data Stewardship Portal**.
3. On the **Sign in** page, enter your Spectrum user name and password to sign in.

Access the Data Stewardship Portal from the Spectrum Home page

To perform this procedure, you need to know the name of the host server running Spectrum Technology Platform and the port number used by Spectrum Technology Platform.

1. Open a web browser and navigate to `http://SpectrumServer:port/`.
For example,
`http://myserver:8080`
2. Click **Spectrum Data Quality > Data Stewardship Portal > Open the Data Stewardship Portal**.
3. On the **Sign in** page, enter a valid user name and password to sign in.

Access the Data Stewardship Portal directly

To perform this procedure, you need to know the name of the host server running Spectrum Technology Platform and the port number used by Spectrum Technology Platform.

1. Open a web browser and navigate to `http://SpectrumServer:port/bsmportal`.
For example,

`http://myserver:8080/bsmportal`

2. On the **Sign in** page, enter a valid user name and password to sign in.

The Dashboard page

The Dashboard page in Data Stewardship Portal displays data that summarizes the status of records belonging to you and other users.

The Dashboard page is displayed by clicking **Dashboard** on the **Data Stewardship Portal** page. The Dashboard displays status and statistics for a selected user. The default is the currently logged in user. You choose to view data for other users if you have modify permissions. A different user may be selected in the box at the top of the page. Only users who currently have exceptions assigned to them appear in users box.

Note: The progress charts for user and dataflow will not appear if you have turned off progress tracking on the Data Stewardship Settings page of Management Console.

Total	Shows total number of records for the currently selected user.
Remaining	Shows number and percentage of records that remaining for the currently selected user.
Resolved	Shows the number and percentage of records that have been resolved for the currently selected user.
Status	The chart shows the selected user's percentages of resolved and remaining records.
Daily Progress	<p>The selected user's progress by day, week, or month. This shows number of records accepted by the user. For an approval flow, it shows the number of records accepted and rejected by the user.</p> <ul style="list-style-type: none"> • Click the Display table button  to view progress in table format. • Click the Display chart button  to view progress in chart format. • Click the Options button  to select the duration of Day, Week, or Month.
Breakdown by dataflow	This section lists active dataflows. To edit a dataflow on the Data Stewardship Portal Editor page, click the dataflow name. To filter the list, enter a portion or all of the characters in a dataflow name in the Filter box  .
Resolved Timeline	This section shows the resolved status of approval flows. Information includes, the ratio and percentage of resolved records, the number of records waiting to be resolved or accepted, the number of records resolved, and the total number of records.

The Data Stewardship Portal Editor

The Data Stewardship Portal **Editor** enables users to manually review, modify, and resolve records that cannot be corrected by stages in a dataflow. Users can also use the Data Stewardship Portal **Editor** to resolve duplicate records and use search tools to look up information that assists in editing, resolving, and rerunning records. A user with modify permission is able to correct incorrect or missing data in a record and resolve a record. Resolved records can then be read from the Data Stewardship repository using the Read Exceptions stage for additional processing.

When an approval flow is defined for data, the Exception Monitor assigns records to users in a sequence of one or more acceptance levels. Exception Monitor initially assigns records to the default user, referred to as the *data steward*. The data steward reviews and edits a record before accepting it. After the data steward accepts a record, it is assigned to the next level of the approval flow. In each level after the data steward, a user can modify, accept, or reject a record. A record that is rejected at any level is assigned back to the data steward. After the last-level user approves and saves a record, it is set to resolved. The resolved records can then be read by the Read Exceptions stage for additional processing and are purged from the Data Stewardship repository.

Customizing Data Stewardship Portal Editor Contents

There are multiple ways you can customize what is shown on the Data Stewardship Portal **Editor** page. You select settings on the side panel to view records in the exception repository for a particular **User** (if the current user has access to other user exception assignments), **Dataflow name**, **Stage label**, **Job ID**, **Status**, **Date**, or **Condition** or **Data Domain** or **Metrics**. You can use the Set filter  tool to only display records that meet certain criteria. The Configure view  options allow you to show or hide fields in the editor.

Do not forget that after you configure settings described in the previous paragraph in combination with the **items per page** setting, there may still be multiple pages of results that are not shown on the initial screen. For example, you set **items per page** to 10 and then apply a field filter to return only the records that have a specific postal code. You may initially see only 10 records, but should you check the bottom right corner, as it may show that there are additional items in the repository.

Search Options

Search options determine which records are visible in the Exception **Editor**.

Click the **SEARCH** tab on the left side of the Data Stewardship Portal **Editor** page to expand or collapse search options. Click any of the following headings on the **SEARCH** tab to expand or collapse criteria for an option.

User	The ID of the user assigned to the dataflow. This information will be visible only if you have modify permissions. This selection is required.
Dataflow name	The name of the dataflow that generated the records. This selection is required.
Stage label	The user-defined name given to the Exception Monitor stage in the dataflow. This information is particularly useful in cases where a dataflow contains multiple Exception Monitor stages. If the person who created the dataflow gave each Exception Monitor stage a meaningful label you can identify which Exception Monitor produced the record. The default label is "Exception Monitor". This selection is required.
Job ID	A numeric identifier assigned to a job by the system. Each time a job runs it is assigned a new job ID. You can choose to view records for all jobs or select a specific job.
Status	Lists possible approval status categories for records in the exception repository. The number of records currently in a category is displayed next to the status label. These values are updated each time a user saves changes. Check the check box next to a status label to display records that currently display the status in the Record status field. You can check multiple check boxes on this panel.
Date	Choose to display records from dataflows run on specified dates. You can specify either or both start and end dates. Click the date button  select a date. Click the time button  to specify the time. The default is 12:00 AM.
Condition	<ul style="list-style-type: none"> • Name—Choose a condition name to only view records that were created as a result of a particular condition. If you specify a name, then both Data Domain and Quality Metrics are set to All. Note: To display condition names for each record, see Search Options—Condition Name. • Data Domain—Specifies the kind of data to which the record belongs, such as Customer or Product. Examples of data domains include Name, Address, and Phone Number. This selection helps you identify which fields in the record require editing. If you specify a data domain, then Name is set to All. • Quality Metrics—The quality measurement that failed the record. Examples of quality metrics include Accuracy, Completeness, and Uniqueness. This information helps you determine why the record was identified as an exception. If you specify a quality metric, then Name is set to All.

Field filters

Field filters control limit records that are displayed in the Data Stewardship Portal **Editor** according to field values.

After records are loaded into Data Stewardship, you filter fields to display only those records in which you are interested. The Data Stewardship Portal initially displays all records from a single Spectrum Technology Platform dataflow. Field filters display only those records that meet specified criteria for

one or more fields. Once filtering is defined for a dataflow, it is applied each time you open the dataflow on the **Editor** page until you clear the filters.

Note: You must apply any selection options described in [Search Options](#) on page 167 before using the field filter tool.

Filters are set by clicking the Set filter button  on the toolbar. Doing so displays the Filter icon  on field headings for fields that can be filtered. You can click the Filter icon  on a field heading to set a filter for the field. After you set a filter on a field, its Filter icon  displays a blue background:

City	FirstName
KEENE	LAKSHMI
KEENE	LAKSHMI

Filters determine which records are visible in both Tabular view and Form view, although they can only be created in the Tabular view. The filter icon on the Form view status bar signifies that records are currently filtered:



Related concepts

[Lookup Tables \(Data Quality Guide\)](#) on page 150

Related tasks

[Setting field filters](#) on page 169

This task describes how to filter the list of records according to field values.

Related reference

[Field filter operators](#) on page 170

This topic describes operators that you may apply in a field filter. The available operators depend on whether a field is a text, date, or numeric field.

Setting field filters

This task describes how to filter the list of records according to field values.

1. On the Data Stewardship Portal **Editor** page, make sure that you are in Tabular view .
2. Click the Set filter button  on the toolbar.
This displays the Filter icon  on every field heading that can be filtered.
3. Click the Filter icon  on a field heading.
4. In the **Show items with value that** list, select the operator that you want to use.
For more information, see [Field filter operators](#) on page 170.
5. Enter one or more values for the filter, then press the Enter key.
If a lookup table was defined in the Data Stewardship settings and is associated with the column, click the value box to expand the list and select values. Otherwise, you can type and add values to the value box.

- Click the **Filter** button to apply the criteria.

Only records whose data matches the filter criteria are then visible. If you add multiple values to a filter, records are displayed whose data matches the criteria for any of the values. All records are hidden if there is no match.

- To filter on additional fields, repeat steps **3** on page 169 through **6** on page 170 for each field that you want to filter.

When you set filters on multiple fields, records are displayed whose data matches the criteria set for all of the fields.

You can click the Set filter button  again to clear filters currently set on fields.

Related concepts

[Lookup Tables \(Data Quality Guide\)](#) on page 150

Related tasks

[Clearing field filters](#) on page 170

You can clear a single field filter or all field filters.

Related reference

[Field filter operators](#) on page 170

This topic describes operators that you may apply in a field filter. The available operators depend on whether a field is a text, date, or numeric field.

Clearing field filters

You can clear a single field filter or all field filters.

- To clear the filter from a single field in Tabular view, first click the Filter icon  on the field heading, then click the **Clear** button.
- To clear filters from all fields in either Tabular view or Form view, click the Set filter button  on the toolbar.

Related tasks

[Setting field filters](#) on page 169

This task describes how to filter the list of records according to field values.

Field filter operators

This topic describes operators that you may apply in a field filter. The available operators depend on whether a field is a text, date, or numeric field.

Is equal to	Looks for records that have exactly the value you specify. This can be a numeric value or a text value. For example, you can search for records with a MatchScore value of exactly 82, or records with a LastName value of "Smith."
Is not equal to	Looks for records that have any value other than the one you specify. This can be a numeric value or a text value. For example, you can search for records with any MatchScore value except 100, or records with any LastName except "Smith."

Starts with	Looks for records that start with a particular value in the selected field. For example, if you filter for "Van" in the LastName field you would see records with "Van Buren", "Vandenburg", or "Van Dyck".
Contains	Looks for records that contain the value you specify in any position within the selected field. For example, if you filter for "South" in the AddressLine1 field, you would see records with "12 South Ave.", "9889 Southport St.", "600 South Shore Dr.", and "4089 5th St. South".
Does not contain	Looks for records that do not contain the value you specify in any position within the selected field. For example, if you filter for "South" in the AddressLine1 field, you would not see records with "12 South Ave.", "9889 Southport St.", "600 South Shore Dr.", and "4089 5th St. South".
Ends with	Looks for records that end with a particular value in the selected field. For example, if you filter for records that end with "burg" in the City field, you would see records with "Gettysburg", "Fredricksburg", and "Blacksburg".
Is greater than	Looks for records that have a numeric value that is greater than the value you specify.
Is greater than or equal to	Looks for records that have a numeric value that is greater than or equal to the value you specify. For example, if you specify 50, you would see records with a value of 50 or greater in the selected field.
Is less than	Looks for records that have a numeric value that is less than the value you specify.
Is less than or equal to	Looks for records that have a numeric value that is less than or equal to the value you specify. For example, if you specify 50, you would see records with a value of 50 or less in the selected field.
is after	Looks for records that have a date or time value that is later than the value you specify.
Is after or equal to than	Looks for records that have a date or time value that is equal to or later than the value you specify.
Is before	Looks for records that have a date or time value that is earlier than the value you specify.
Is before or equal to than	Looks for records that have a date or time value that is equal to or earlier than the value you specify.

Related tasks

[Setting field filters](#) on page 169

This task describes how to filter the list of records according to field values.

Viewing Records

To open the Data Stewardship Portal **Editor** to view records, click **Editor** on the Data Stewardship Portal menu bar. The Data Stewardship Portal **Editor** will also open when you click a dataflow in the **Breakdown by dataflow** list on the Dashboard.

On the **Editor** page, you can view records in either Tabular view or Form view.

Tabular view

In Tabular view you can load up to 100 records per page. You can change how many records are displayed per page in the **items per page** box at the bottom of the screen. In this view, you can scroll through the list and edit records in any order. When you finish editing records on a page, you can click the Save changes button  to simultaneously save changes to all of the records on a page.

Form view

In Form view you view and edit data for one record at a time. You can click the navigation buttons to scroll between records in the current dataflow. In this view you must save the edits to a record before you proceed to a different record. You cannot save edits more than one record at a time in Form view.

Viewing Record Details

Regardless whether you use Tabular view or Form view, the Exception grid shows all the fields for a record as well as record status, exception type, and any comments that have been added to the record. In the Tabular view, you can view additional details about a record by clicking the arrow in the left column of a record. In the Form View, click **Show detail** beneath the form fields. In either view, expanding details for a record shows the following tabs.

Details tab on page 172	This tab shows details about a record.
Condition tab on page 173	The tab shows information about the condition that identified the record as an exception.
History tab on page 173	This tab shows historical information for a record.
Approval flow tab on page 174	This tab shows the current status of a record and the levels that have been completed in an approval flow.

Details tab

This tab shows details about a record.

Job ID	A numeric identifier assigned to a job by the system. Each time a job runs it is assigned a new job ID.
Group by	If the dataflow was configured to return all records in the records group, this shows the field by which the records are grouped. This only applies to dataflows that perform matching, such as dataflows that identify duplicate records or dataflows that group records into households.
Created on	The date and time when the Exception Monitor identified the record as an exception.

Record type	The type of record. This is one of the following: <ul style="list-style-type: none"> • Exception • Group exception • Group non-exception • Best of Breed
Resolved	Shows whether record is resolved (Yes or No).
Approval flow types	Shows inclusion in approval flow types.

Condition tab

The tab shows information about the condition that identified the record as an exception.

Name	The name of the condition that identified the record as an exception. Condition names are defined by the person who set up the dataflow.
Domain	The kind of data that resulted in an exception. Examples of data domains include Name, Address, and Phone Number. This information helps you identify which fields in the record require editing.
Metrics	The quality measurement that the record failed. Examples of quality metrics include Accuracy, Completeness, and Uniqueness. This information helps you determine why the record was identified as an exception.

History tab

This tab shows historical information for a record.

Every time a record is changed, field values and historical information are retained by the system. The table on this tab shows saved versions of a record. The column heading for a field is highlighted in blue when its value was changed at any time in the history of the record. In these highlighted columns, a field value is highlighted in blue when it does not match the value shown in the preceding version of the record.

The initial columns in the table provide the following information about each version of a record:

Version	The revision number of the change. This starts with 1.0 for the original version of a record and increments 1.1, 1.2, 1.3, and so forth each time that changes are saved to a record.
Modified time	Date and time when record was last saved.
Modified by	The user who made the last change.
Assigned to	The user to whom the record was assigned at the time of its revision.
Record status	The current status of the record. This is one of the following: <ul style="list-style-type: none"> • Accepted—Record has been accepted in an approval flow. • Edited—Record has been edited. • Failed—Record failed revalidation. • New—Record has not been reviewed or edited.

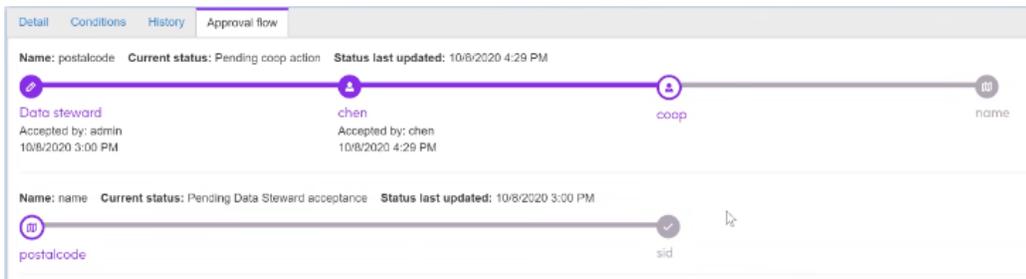
- Pending—Record has been reviewed or edited but has been retained for additional review.
- Rejected—Record has been rejected in an approval flow.
- Resolved—Accepted record has been saved for additional processing.

Condition name	The condition in the Exception Monitor stage that generated the record.
Comment	The most recent comments that were entered for the record.

Approval flow tab

This tab shows the current status of a record and the levels that have been completed in an approval flow.

Shows the following information for an approval flow. When the record is in multiple approval flows, the information and visual representation is provided for each flow. When a record is flagged by more than one condition, it progresses through each approval flow. Each status line represents a single approval flow plan.



Name	The name of the condition.
Current status	Shows the whether a record is pending action by a reviewer, rejected, accepted, or resolved.
Status last updated	Shows the date and time that the status was last updated.

Sorting records

In Tabular view, you can sort records by values in a single column. You can sort records by text (A to Z or Z to A), numbers (smallest to largest or largest to smallest), and dates and times (oldest to newest and newest to oldest).

1. In Tabular view, choose the field that you want to use to sort records.
2. Click the column heading to sort records based on ascending order of values in the column. This displays the ascending sort ▲ icon next to the column heading to indicate that records are sorted based on ascending values in the column.
3. To sort based on a descending sort of values in the column, click the column heading again. This displays the descending sort ▼ icon next to the column heading to indicate that records are sorted based on descending values in the column.

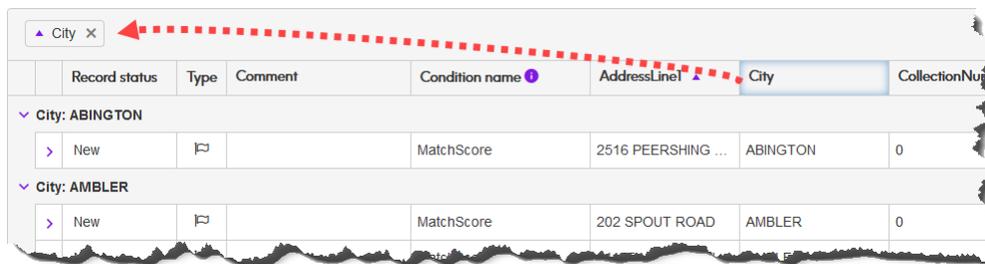
You can sort on any column except the **Condition name** column, which is displayed in the record table when the **Condition name** field check box is checked in the Configure view. For more information, see [Configuring the view](#) on page 175.

Grouping records

In Tabular view you can group records by values in one or more columns. When you group by values in more than one column they are grouped hierarchically in the order that you choose.

To complete this task, the **Allow grouping records list** check box must be checked in the **Configure view** options. For more information, see [Configuring the view](#) on page 175.

1. For the first column on which to group records, drag the column heading to the panel above the column headings.



Completing this step groups records by values within this column.

2. Optional: To change the sort order of values in the group, click the ascending ▲ or descending ▼ icon on the column entry in the panel.

This changes the order for records displayed on a page.

Note: To group records and maintain sort order across pages, sort the column before adding it to the grouping. For more information, see [Sorting records](#) on page 174.

3. To add subordinate groupings, you can drag additional column headings to the panel above the table.

Completing this step creates groupings in the order that you add headings.

4. Optional: To change the order in which records are grouped, drag a column heading already in the panel before or after another column heading in the panel.

Configuring the view

You can select which fields appear and change the order in which they appear by clicking the Configure view button . This displays the **Configure view** options.

Changes to the **Configure view** settings are saved for a user and dataflow name, so that when you reopen a dataflow the configuration from the previous session is applied. Changes made here apply to editing records in both Tabular view and Form view.

Filter

Enter all or part of a field name in the Filter box  to only display field names that contain the entered characters in their names. The search is case insensitive and the text can be anywhere in the name of a field.

All	List all fields. Check or uncheck a field to expose or hide it in both Tabular view and Form view.
Visible	Lists fields that are currently checked in the Field name column.
Hidden	Lists fields that are currently not checked in the Field name column.
Display Condition name field	<p>In the Field name column, check a check box to display a column in both Tabular view and Form view. This field displays the name of the condition that caused the record to be added to the repository. Records that meet criteria for multiple conditions will display condition names as a comma separated list.</p> <p>Note: To only display records for a single condition, see the Search Options—Condition Name.</p>
Change field order	<p>The fields are listed in the Configure view dialog box in the same order that they are displayed in the records table. You can change the order in which fields appear by dragging and dropping fields in the desired order. However, you cannot rearrange fields while there is text in the Filter box . You must first clear the contents of the Filter box before you can reorder fields.</p> <p>Note: You can also order fields in Tabular view. Drag column headings on the record table to reorder the columns. Any changes to the table itself will also appear in the Configure view dialog box the next time that you open it.</p>
Configure field widths	<p>You can specify the width of columns that appear in the Form view by entering values in the Form columns boxes. Values for adjacent visible fields in this column must sum to 12. For example, if you wanted four fields to be placed side-by-side in four equal-length fields on a single row, you could give each field a column width of "3." Likewise, if you wanted to include just two fields of equal length in a row, you could enter "6" for both fields under Form columns. These proportional settings are illustrated in the first figure, Figure 1: Form columns configured for a rowupdate image, checkboxes changed, first option changed to "Optimize record editors for keyboard navigation" on page 177. The fields in the figure are set to 4, 3, 2, and 3 to define four adjacent fields of various widths that sum to 12. The second figure, Figure 2: Fields in Form view rendered as configured on page 177, shows the result from Form view.</p>

Configure view ×

▼

All
Visible
Hidden

Field name	Form columns
↕ <input checked="" type="checkbox"/> Approved	Default ▼
↕ <input checked="" type="checkbox"/> Type	Default ▼
↕ <input checked="" type="checkbox"/> Comment	Default ▼
↕ <input checked="" type="checkbox"/> AddressLine1	4 ▼
↕ <input checked="" type="checkbox"/> City	3 ▼
↕ <input checked="" type="checkbox"/> State	2 ▼
↕ <input checked="" type="checkbox"/> PostalCode	3 ▼

Optimize tabular editor for keyboard navigation
 Allow grouping records list

OK
Cancel
Reset to default

Figure 1: Form columns configured for a row

Approved
 Comment

AddressLine1
555.55HC71 BOX60

City
CTR STRAFFORD

State
NH

PostalCode

Figure 2: Fields in Form view rendered as configured

Note: When you configure field widths under **Format columns**, the configured fields will no longer be responsive to the size of the browser window. In contrast, fields with the **Default** setting are always responsive to the size of browser window, wrapping and unwrapping when the window is resized.

Enable Keyboard Navigation

Check the **Optimize tabular editor for keyboard navigation** check box to use the Tab key to navigate records in the Tabular view of the Data Stewardship Portal **Editor** page.

Group records in Tabular view

Check the **Allow grouping records list** check box to group records by their column values. This displays a panel above the column headings in which to configure groupings. Records may be grouped hierarchically in descending or ascending order on a page by values in one or more columns. This check box is checked by default. For more information, see [Grouping records](#) on page 175.

Reset to Default

Click the **Reset to default** button to check all of the check boxes in the **Field name** column and to set **Form columns** to "Default".

Editing and resolving records

The purpose of editing a record is to correct or add information and then accept and resolve the record so that it can be processed successfully. Editing a record may involve using other Spectrum Technology Platform services or consulting external resources such as maps, the Internet, or other information systems in your company.

After reviewing records, you can edit, accept, and resolve them directly on the Data Stewardship Portal **Editor** page. You can edit one or more records simultaneously in the Tabular view, and you can edit one record at a time in the Form View.

Note that read-only fields cannot be edited. If you want to make a read-only field editable, you would need to delete all records for that dataflow and job ID and run the dataflow again after configuring the fields accordingly in the Write Exceptions stage. This would produce new records with editable fields. Furthermore, you cannot edit a field to contain a value that does not match its data type. For example, you cannot enter alphabetic characters in a field that has a numeric data type.

Editing records in Tabular view

Perform steps in this task to edit and approve records in Tabular view.

1. In Data Stewardship Portal, open records for the dataflow that you want to edit on the **Editor** page.
To do this, click **Editor** on the Data Stewardship Portal home page, then in the **Dataflow name** box on the sidebar, select the name of the dataflow.
2. Make sure that the records are displayed in Tabular view.
Click the Tabular view button  on the **Editor** toolbar.
3. You can now edit field values for a single record or for multiple records:
 - [To edit a field for a single record in Tabular view](#) on page 178
 - [To edit a field for multiple records in Tabular view](#) on page 179

To edit a field for a single record in Tabular view

1. In any record, click to edit a value.
Right-click the field to access **Cut**, **Copy**, and **Paste** commands. A field will be grayed out if it cannot be edited.

To undo a change, select the record you want to undo and click the Undo changes button  on the toolbar.

After you edit a field, green triangle  in the upper-left corner that indicates edits are not yet saved. After you save changes to the record, the status changes to Edited.

2. Optional: You can add comments about your changes in the **Comment** box.
Comments are visible to other users and can be used to help keep track of the changes made to the record or to recommend additional changes.
3. After you are satisfied with corrections made to a record, click Accept .
This changes the **Record status** to Accept.

Note: Anytime before you click the Save button, you can click Accept  or Reject  to change the status from Accept to Pending. This allows you to continue editing the record at a later time.

4. Click the Save changes button  when you are finished editing records.
This changes status for the record to Resolved. The records may not appear in the table if they are no longer available or no longer match the search/filter criteria.

To edit a field for multiple records in Tabular view

1. Click the field in the first record that you want to change.
2. Press Shift to extend the selection to multiple records or Ctrl to add a single record to the selection, and click the same field in another record.
For instance, if you want to change all instances of "L.A." to say "Los Angeles," click within the City field for one of the records. Then press Ctrl or Shift and click within that same column to extend or add to the selection.
3. Edit the field value in the last record that you selected.
Right-click the field to access **Cut**, **Copy**, and **Paste** commands. A field will be grayed out if it cannot be edited.

To undo a change, select the records you want to undo and click the Undo changes button  on the toolbar.

After you edit a field, green triangle  in the upper-left corner that indicates edits are not yet saved. After you save changes to the record, the status changes to Edited.

4. Optional: You can add comments about your changes in the **Comment** box.
Comments are visible to other users and can be used to help keep track of the changes made to the record or to recommend additional changes.
5. After you are satisfied with corrections made to a record, click Accept .
This changes the **Record status** to Accept.

Note: Anytime before you click the Save button, you can click Accept  or Reject  to change the status from Accept to Pending. This allows you to continue editing the record at a later time.

- Click the Save changes button  when you finish editing records on a page. This changes status for the records to Resolved. The records may not appear in the table if they are no longer available or no longer match the search/filter criteria.

Editing fields in Form view

Perform steps in this task to edit records in the Form view.

- In the Data Stewardship Portal **Editor** page, display records in the exception repository for the dataflow that you want to edit .
To do this, click **Editor** on the Data Stewardship Portal toolbar, then in the **Dataflow name** box on the sidebar select the name of the dataflow.
- Make sure that the records are displayed in Form view.
Click the Form view button  to display fields in the first record.
- In any field, click to edit a value.
Right-click a field to access **Cut**, **Copy**, and **Paste** commands. A field will be grayed out if it cannot be edited.

To undo a change, select the record you want to undo and click the Undo changes button  on the toolbar.
An edited field is highlighted green to indicate unsaved changes. After you save changes to the record, the status changes to Edited.
- Optional: You can add comments about your changes in the **Comment** box.
Comments are visible to other users and can be used to help keep track of the changes made to the record or to recommend additional changes.
- After you are satisfied with corrections made to a record, click Accept .
This changes the **Record status** to Accept.

Note: Anytime before you click the Save button, you can click Accept  or Reject  to change the status from Accept to Pending. This allows you to continue editing the record at a later time.

- Click the Save button  .
This changes status for the records to Resolved. You will either see the same record containing your changes or you will see the next record in the list because the record you changed is no longer available or no longer matches the search/filter criteria.
- Use the navigation buttons at the bottom of the form to go to previous or next record.
You can also go directly to the first or last record.

Exporting Records

You can export data for all records from a dataflow to a comma separated value (CSV) file. You can then import the data into other applications such as Microsoft Excel.

This procedure saves data for all of the records from a dataflow to a CSV file. The first (or header) row in the CSV file lists the field names. Subsequent rows list the corresponding values for each record. Field names and data values are enclosed between double quotation marks ("").

1. In Data Stewardship Portal, open records for the dataflow that you want to export on the Data Stewardship Portal **Editor** page.

To do this, click **Editor** on the Data Stewardship Portal toolbar, then in the **Dataflow name** box in the sidebar select the name of the dataflow.

2. Make sure that the records are displayed in Tabular view.

If necessary, click the Tabular view button  on the **Editor** toolbar.

3. On the **Editor** toolbar, click the Export all records button .

4. If prompted, choose to save the file.

All the records created by the dataflow are exported to the browser default download folder. The exported file name is *DataflowName_StageLabel.csv*, with the dataflow name and stage label of the dataflow that created the records. If records created by a dataflow are exported repeatedly, the files are numbered incrementally in parentheses after the dataflow name as

DataflowName_StageLabel(number).csv.

After you complete this procedure, you can open the file and edit records in another application, such as Microsoft Excel or a text editor. After you complete the changes, you can load the data back into a workflow in Enterprise Designer or Flow Designer.

Approval flows

Approval flows assign records to sequential levels for review and acceptance by a user.

Records are initially assigned to the data steward specified in the Exception Monitor stage. The data steward can review, edit, and accept records. If a record is accepted, it is then assigned to the first approver as defined in the approval flow plan. At this level, the approver can Accept or Reject a record. If a record is accepted it is be assigned to the next approver in the approval flow. If a record is rejected, it is assigned back to the data steward. This is repeated at subsequent levels until the final approver accepts and resolves the record. Resolved records can then be read using the Read Exception stage for additional processing.

Creating approval flows

This process defines tasks that must be completed to create an approval flow. The approval flow process is illustrated in the [Approval Flow Process diagram](#).

- [Define approval flow types](#) on page 182

Define approval flow types in the Data Stewardship Settings of the Management Console.

- **Specify the approval flow type for a condition** on page 182

Assign an approval flow type to a condition to be used in an approval flow. This task is performed in the Exception Monitor stage. Perform this procedure in the Exception Monitor settings in an Enterprise Designer dataflow.

- **Enable approval flow types in security roles** on page 183

Execute permission must be applied to an approval flow type within a security role before you can enable it in an approval flow. Perform this procedure in the Roles settings of the Management Console.

- **Assign approval flow role to new or existing users** on page 183

To define reviewers in an approval flow, assign an approval flow role to new or existing users. Perform this procedure in the Users settings of the Management Console.

- **Define an approval flow** on page 184

Define levels and specify the role and user for each level in the Management Console Approval Flow settings. Perform this procedure in the Approval Flows settings of the Management Console.

Define approval flow types

Define approval flow types in the Data Stewardship Settings of the Management Console.

1. Open the Spectrum Technology Platform **Management Console**.
On the Spectrum Technology Platform **Home** page, click **Platform Client tools > Web > Open Management Console**.
2. On the **Management Console** page, click **Resources > Data Stewardship Settings**.
3. On the **Data Stewardship Settings** page, click **Approval Flow Types**.
4. Click the Add type button to create a new approval flow type.
5. In the **Name** box, enter a meaningful name.
6. Click **Save**.

Specify the approval flow type for a condition

Assign an approval flow type to a condition to be used in an approval flow. This task is performed in the Exception Monitor stage. Perform this procedure in the Exception Monitor settings in an Enterprise Designer dataflow.

1. In Enterprise Designer, double-click the Exception Monitor stage in a workflow to open its properties.
2. On the **Conditions** tab, click a condition and click the **Modify** button.
3. Under **Approval flow**, in the **Type** drop-down box, click an approval flow type.
4. Optional: To specify the data steward for the approval flow, click to choose a user in the **Assign to** box.

The default data steward for an approval flow, if you leave the **Assign to** setting blank, is the user who runs a job.

- Click **OK**, and click **OK** again.
- Click **File** > **Save** to save the changes to the workflow.

Enable approval flow types in security roles

Execute permission must be applied to an approval flow type within a security role before you can enable it in an approval flow. Perform this procedure in the Roles settings of the Management Console.

- Open the Spectrum Technology Platform **Management Console**.
On the Spectrum Technology Platform **Home** page, click **Platform Client tools** > **Web** > **Open Management Console**.
- On the **Management Console** page, click **System** > **Security**.
- On the **Security** page, click the **Roles** tab.
- Click the Add role button **+**.
- Click to expand the **Data Stewardship** section.
Approval flow types defined in the **Data Stewardship Settings** show in the table.
- For approval flow types that you want to include in approval flows, check the check box in the **Execute** column.
Execute permission is required to create an approval flow.

Assign approval flow role to new or existing users

To define reviewers in an approval flow, assign an approval flow role to new or existing users. Perform this procedure in the Users settings of the Management Console.

- Open the Spectrum Technology Platform **Management Console**.
On the Spectrum Technology Platform **Home** page, click **Platform Client tools** > **Web** > **Open Management Console**.
- On the **Management Console** page, click **System** > **Security**.
- On the **Security** page, click the **Users** tab.
- Open settings for a new or existing user.

To	Do this
Add a new user	Click the Add user + button to open the Add User page.
Edit an existing user	Click the check box next to a user and click the Edit user button  to open the Edit User page.

- For a new user, enter the user name and password.
- In the **Roles** box, check the check box next to the approval flow role that you want to use (see [Enable approval flow types in security roles](#) on page 183).
- Make sure that **Enabled** is set to **ON**.

8. Click the **Save** button.
9. Repeat steps **4** on page 183 through **8** on page 184 for each user that you want to include in the approval flow.

Define an approval flow

Define levels and specify the role and user for each level in the Management Console Approval Flow settings. Perform this procedure in the Approval Flows settings of the Management Console.

1. In the **Management Console**, click **Resources > Approval Flows**.
Approval flow types are listed in the **Type** column of the **Approval Flows** table.
2. In the **Type** column, click an approval flow type that you want to configure as an approval flow. This opens the Edit Approval flow page for the selected approval flow type.
3. Click the Add approval level button .
This adds a level to the approval flow. Each level will be associated with a reviewer defined in .
4. In the **Roles** column, select the role to be used for the level.
5. In the **User** column, select the user that records will be assigned to for this level.
If no user is selected, the record is assigned to the first user in alphabetical order.
6. In the **Days to Approve** column, select the maximum number of days in which the user must review and edit records.
7. To add additional levels to the approval flow, repeat steps **3** on page 184 through **7** on page 184.
8. Click to switch the **Enabled** box to **ON**.
9. Click the **Save** button.

Reviewing records in an approval flow

Perform steps in this task to edit, accept, or reject records in Tabular view.

1. On the **Data Stewardship Portal** Dashboard page, click the dataflow to view records that you want to review.
Alternatively, click **Editor**, then in the **Dataflow name** menu on the sidebar select the name of the dataflow.
2. Make sure that the records are displayed in Tabular view.
Click the Tabular view button  on the Editor toolbar.
3. In any record, click to edit a value.
Right-click the field to access **Cut**, **Copy**, and **Paste** commands. A field will be grayed out if it cannot be edited.
If you need to undo changes you made, select the record you want to undo and click the Undo changes  button.

After you have edited a field, you will notice a green triangle  appear in the upper-left corner of the field. This is a visual cue to remind you that the value of the field has been changed but is not yet saved.

4. Optional: You can add comments about your changes in the **Comment** box.
Comments are visible to other users and can be used to help keep track of the changes made to the record or to recommend additional changes. In an approval flow, you can describe additional changes necessary for approval.
5. Click the Save changes button .
This changes the record status to Edited.
6. After you have reviewed or edited values, you can click to Accept  or Reject  a record.
Clicking Accept  changes the record status to Accepted. Clicking Reject  changes the record status to Rejected.

Note: Anytime before you click the Save button , you can click Accept  to change the status from Accepted to Pending. This allows you to continue editing the record at a later time. Clicking Reject  changes the status to **Edited**

7. Click the Save changes button  button when you are finished editing records on a page.
Accepted records are assigned to the next approver in the approval flow. For the last approver in the approval flow, the record status is changed from Accepted to Resolved, and the records are made available for additional processing using the Read Exceptions stage.
Rejected records are assigned back to the data steward.

Resolving Duplicate Records

Duplicate resolution exceptions occur when the Spectrum Technology Platform cannot confidently determine whether a record is a duplicate of another. There are three ways to resolve duplicate records.

Note: Duplicate records can only be resolved with the Resolve Duplicates function on the Tabular view. However, you can still edit those records in the Form view.

Note: If you are working with Duplicate records in an Approval flow, see the section on [Approval Flows](#).

One approach is to group duplicate records together into collections. When you approve the records they can then be processed through a consolidation process to eliminate the duplicate records in each collection from your data.

Another approach is to edit the records so that they are more likely to be recognized as duplicates, for example correcting the spelling of a street name. When you approve the records, Spectrum Technology Platform reprocesses the records through a matching and consolidation process. If you corrected the records successfully, Spectrum Technology Platform will be able to identify the record as a duplicate.

Yet another approach to resolving duplicate records is to create a best of breed record. This combines the other two approaches by managing record collections and then editing one of the records in the collection to include fields from both the original and duplicate records. This record is then known as the best of breed record.

Making a Record a Duplicate of Another

Duplicate records are shown as groups of records in the Data Stewardship Portal. You can make a record a duplicate of another by moving it into the same group as the duplicate record.

To make a record a duplicate:

1. In the Exception Editor, select the record you want to work on then click the Resolve duplicates button .

The **Duplicate Resolution** view shows duplicate records. In this view, the initial record is highlighted in bold text. A column label is highlighted in blue text on any column of user data that contains values different from the value in the record. In these columns, values that do not exactly match the record are also highlighted in blue. Keep in mind that a mismatch may not be visible when a value is null or consists only of spaces.

For each record, the **MatchRecordType** column specifies one of three types of exceptions:

Suspect	A record that other records are compared to in order to determine if they are duplicates of each other. Each collection has one and only one suspect record.
Duplicate	A record that is a duplicate of the suspect record.
Unique	A record that has no duplicates.

2. If necessary, correct individual records as needed.
For more information, see [Editing and resolving records](#) on page 178.
3. In the **CollectionNumber** or **CandidateGroup** field, enter the number of the group that you want to move the record into. Alternatively, you can drag and drop records between groups.

The record is made a duplicate of the other records in the group.

In some cases you cannot move a record with a MatchRecordType value of "suspect" into another collection of duplicates.

Note: Records are grouped by either the CollectionNumber field or the CandidateGroup field depending the type of matching logic used in the dataflow that produced the

exceptions. Contact your Spectrum Technology Platform administrator if you would like additional information about matching.

4. When you are done modifying the record, click Accept . This changes the record status to Accept. If you click Accept  again, the status is changed to Pending, so that the record can be reviewed later.
5. To save your changes, click the Save changes  button. This changes the record status to Resolved and makes the record available for additional processing using the Read Exceptions stage.

Creating a New Group of Duplicate Records

In some situations you can create a new group of records that you want to make duplicates of each other. In other situations you cannot create new groups. Your ability to create new groups is determined by the type of Spectrum Technology Platform processing that generated the records.

1. In the Editor, select the record you want to work on then click the Resolve duplicates button .

The **Duplicate Resolution** view shows duplicate records. In this view, the initial record is highlighted in bold text. A column label is highlighted in blue text on any column of user data that contains values different from the value in the record. In these columns, values that do not exactly match the record are also highlighted in blue. Keep in mind that a mismatch may not be visible when a value is null or consists only of spaces.

For each record, the **MatchRecordType** column specifies one of three types of exceptions:

Suspect	A record that other records are compared to in order to determine if they are duplicates of each other. Each collection has one and only one suspect record.
Duplicate	A record that is a duplicate of the suspect record.
Unique	A record that has no duplicates.

2. If necessary, correct individual records as needed. For more information, see [Editing and resolving records](#) on page 178.
3. Select a record that you want to put in the new collection then click the Create new collection button .

The new collection is automatically given a unique collection number, and the record you selected becomes a suspect record.

Note: If you do not see the **New Collection** button, you cannot create a new collection for the records you are working with. You can only create new collections if the dataflow that produced the exceptions contained an Interflow Match or an Intraflow Match stage, but not if it contained a Transactional Match stage. Contact your Spectrum Technology Platform administrator if you would like additional information about these matching stages.

4. Place additional records in the collection by entering the new collection's number in the record's `CollectionNumber` field.
5. When you are done modifying the record, click Accept .
This changes the record status to Accept. If you click Accept  again, the status is changed to Pending, so that the record can be reviewed later.
6. To save your changes, click the Save changes  button.
This changes the record status to Resolved and makes the record available for additional processing using the Read Exceptions stage.

Making a Record Unique

Follow steps in this procedure to change a record from a duplicate to a unique.

1. In the **MatchRecordType** field, click "Unique".
2. When you are done modifying the record, click Accept .
This changes the record status to Accept. If you click Accept  again, the status is changed to Pending, so that the record can be reviewed later.
3. To save your changes, click the Save changes  button.
This changes the record status to Resolved and makes the record available for additional processing using the Read Exceptions stage.

Fields Automatically Adjusted During Duplicate Resolution

When you modify records in the duplicate resolution view of the Data Stewardship Portal, some fields are automatically adjusted to reflect the record's new disposition.

Table 8: Records Processed by Interflow or Intraflow Match

Action	Values Automatically Applied to Fields
Moving a record from one collection to another	<p>If you move a record into a collection of duplicates:</p> <ul style="list-style-type: none"> • MatchRecordType: Duplicate • MatchScore: 100 • HasDuplicates: D (This field is only present if the dataflow contained an Interflow Match stage.) <p>If you move a duplicate record into the collection of unique records (collection 0):</p> <ul style="list-style-type: none"> • MatchRecordType: Unique • MatchScore: No change • HasDuplicates: U (This field is only present if the dataflow contained an Interflow Match stage.) <p>If you move a suspect record into the collection of unique records (collection 0):</p> <ul style="list-style-type: none"> • MatchRecordType: Unique • MatchScore: 0 • HasDuplicates: N (This field is only present if the dataflow contained an Interflow Match stage.)
Creating a new collection	<ul style="list-style-type: none"> • MatchRecordType: Suspect • MatchScore: No value • HasDuplicates: Y (This field is only present if the dataflow contained an Interflow Match stage.) <p>Note: If the record came from a dataflow that contained an Interflow Match stage only records with a value of "input_port_0" in the InterflowSourceType field can be a suspect record.</p>

Table 9: Records Processed by Transactional Match

Action	Values Automatically Applied to Fields
Change MatchRecordType to Duplicate	<ul style="list-style-type: none"> • HasDuplicates: D • MatchScore: 100
Change MatchRecordType to Unique	<ul style="list-style-type: none"> • HasDuplicates: U • MatchScore: unchanged

Action	Values Automatically Applied to Fields
Change HasDuplicates to D	<ul style="list-style-type: none"> • MatchRecordType: Duplicate • MatchScore: 100
Change HasDuplicates to U	<ul style="list-style-type: none"> • MatchRecordType: Unique • MatchScore: unchanged
Change HasDuplicates to Y	<ul style="list-style-type: none"> • MatchRecordType: Suspect • MatchScore: blank
Change HasDuplicates to N	<ul style="list-style-type: none"> • MatchRecordType: Suspect • MatchScore: blank

Creating a Best of Breed Record

Follow steps in this procedure to create a best of breed record in a collection.

Best of breed records can only be added to non-zero collections. Also, a best of breed record must be created for each collection within a group before you can accept all the records within that group. If a Best of Breed record does not exist for a collection and the user clicks save, the record status changes to Failed.

1. Select a record in the collection and click the Add Best of Breed Record button ).

This will make a copy of the record you selected and place it just after that record in the Data Stewardship Portal **Editor**. When you save the best of breed record, it will move to the bottom of the record collection. If you look at the best of breed record's history, you will see a link to the record from which it was copied.

Note: You can undo this step by selecting the best of breed record and then clicking the Remove Best of Breed Record button .

2. To combine fields into the best of breed record, drag and drop fields from the duplicate records into the best of breed record.

There are some limitations to drag-and-drop in this environment:

- The cell you are dragging from one record to another cannot be in an "edit" state (though the cell where you are dropping it can be).
- When you drop one cell onto another cell (whether in the same row or a different row) the second cell will take on the value of the first cell.
- You cannot drag and drop Boolean-type cells.

- You can drag non-string-type cells to string-type cells (such as date fields), but you cannot drop string-type cells to non-string-type cells.
- You cannot drop a cell onto another cell that uses lookup values (such as combination box editors), but you can drag a value from these cells to other string-type cells.
- You cannot drop a cell onto a read-only cell.

3. Click the Save changes  button.

Related reference

See [Write Exceptions options for more information on best of breed functionality](#). on page 324

Using Search Tools

The Data Stewardship Portal **Editor** provides search tools to assist you in looking up information that may help you edit and approve records and rerun them successfully. The tools include Maps, OpenNameParser, and other services you have licensed in Spectrum Technology Platform. Tools with icons that have a plus sign in the middle are premium tools and will incur charges when used.

Note: These search tools must first be configured on the **Data Stewardship Settings** page of Management Console. If none are selected on that page, none will appear on this page. Also, regardless of which tools are selected in Management Console, the services that appear to the user here are limited to those for which he has view permissions.

Using Maps

Follow steps in this procedure to search for map locations that correspond to records on the Data Stewardship Portal **Editor** page.

Before you can use this tool, the **Maps** service must be selected on the **Data Stewardship Settings** page of Management Console.

The Maps search tool displays the location of an address on a map. You can move the coordinates to different addresses, zoom in or out, and view adjacent regions.

Note: The Maps search tool is provided by Precisely LI GeoMap and GeoCode API services. You must be connected to the Internet to use this tool.

1. On the Data Stewardship Portal **Editor** page, click the record you want to research.
2. Below the records table, click **Search Tools** to expand the view.
3. In the **Service** field, select **Maps**.
4. To search for a location, do either of the following:
 - Enter an address in the **Input** box.
 - Select check boxes next to the **Field Name** column to search using field values.

For example, to search for an address on a map, you would select AddressLine1 and City. If you only want to view a city on a map, you would select City and StateProvince.

5. Click search button .

The result is displayed in the on the map in the form of a pin . Click the Zoom in  and Zoom out  buttons to adjust the detail. Drag the map to view adjacent regions. The address and the longitude and latitude are displayed in the address box above the map.

10 DOWNING STREET LONDON SW1A 2AA (longitude: -0.12767, latitude: 51.50354)



6. Drag the pin  to view the address and coordinates for another location on the map.

After completing an initial map search, you can click another record and the search button , and the map will update accordingly.

Using Spectrum Service Search Tools

To perform this procedure, you must first configure search tools on the **Data Stewardship Settings** page of Management Console. Otherwise **Search Tools** will not appear on the Data Stewardship Portal **Editor** page.

Precisely service search tools include all services for which you are licensed, such as ValidateAddress, GetPostalCodes, and so on. You can use these services within the Data Stewardship Portal **Editor** to look up and validate exception data that you are attempting to correct or augment.

Note that when using this feature you will only see services if you have view permissions for Services under the Platform group for role security. Likewise, in order to run services, you will need execute permissions for Services. However, these permissions can be modified by using Secured Entity Overrides. Using a combination of top-level permissions and overrides the administrator can manage the list of services that a particular user or role has access to in the Data Stewardship Portal **Service** box.

1. Select the record that contains the data you want to look up.
2. Below the records table, click to expand **Search Tools**.
3. In the **Service** box, choose the service you want to use, such as ValidateAddress or GetCandidateAddresses.
4. If the exception record contains fields used in that service but under different names, map the Service Input fields to the Exception Fields on the **Service Fields** tab.
For example, if you are using ValidateAddress and a record does not include an AddressLine1 field but instead includes an AddrLine1 field, select "AddrLine1" in the Exception Field column of the AddressLine1 row. You must have at least one input field mapped before running the service.

Note: The Data Stewardship Portal remembers the maps you create from service fields to exception fields as long as you are mapping records with the same field names, stage names, and dataflow names. For instance, if a record has a field named "AddrLine1" and you map it to "AddressLine1", it will remember this map as long as you are working with

records that contain "AddrLine1" and that were created in the same stage by the same dataflow.

5. Optional: Repeat step 4 on page 192 for Service Output fields.
This step is optional, but you must have at least one output field mapped before you can apply the service data. Note that you cannot map a read-only Exception Field to the Service Output fields.
6. Click the **Options** tab to select database resources and to view and change service options that were set in Management Console.
If you don't know the purpose of a particular option, hover over that option to see its description. Changes you make here will persist when used by the same user, dataflow, and stage of the record.
Note: If the service you are using requires a database, you must have first configured the database resource in Management Console. For example, if you are reviewing U.S. records using Validate Address, you must have configured a U.S. database in Management Console.
7. Click **Run Service**.
The updated record will appear on the **Result** tab, along with a status code indicating the success of the record. Fields that are mapped to records will be designated by an asterisk.
8. Select the result record and click **Apply Service Data** to transfer the data to the mapped fields of the record.
9. If you want to reprocess the updated record, click the **Approved** check box for that record and then click **Save**.

The Manage page

The **Manage** page on the Data Stewardship Portal enables a user with view and modify permissions to review and manage record activity for all assignees. It also provides the ability to reassign records from one user to another. If you have delete permissions, you can delete an entire group of records from the system based on dataflow name and job ID.

Reviewing Record Activity

In the Data Stewardship Portal, click **Manage** to display the **Manage** page.

Click the **Status** button to view the **Status** section of the **Manage** page. This section shows record activity by assignment or dataflow name. You can specify which one displays by clicking the

Assignments or **Dataflows** button near the top-right corner. **Assignments** provides the number of records assigned to each user as well as how many of those records have been approved and how many remain. It also shows the number of normally processed records that exist within groups that also contain records that have exceptions. Note that non-exception data will appear only if there are normally processed records within a group also containing records with exceptions and if there are no security overrides preventing this data from being shown.

Dataflows provides the number and percentage of records that have been approved for each dataflow. For dataflows containing more than one Exception Monitor stage, you can further break down this information by those stages. For example, if a dataflow contains two Exception Monitor stages, "NameExceptions" might be responsible for 37% of all exceptions in the dataflow and "DataExceptions" might be responsible for the other 63%. In this situation, the status of each **dataflow** is still calculated as a whole as it relates to all dataflows in the repository, while the status of each **stage** is calculated based on the total number of exceptions for the selected dataflow.

Additionally, you can filter the information that displays by entering search criteria in the Filter box. The list will dynamically auto-populate with dataflows or assignees that match the letters you type.

Click the **Progress** button to display the **Progress** section. This section shows the cumulative progress for all users or all dataflows whose individual statuses are shown on the Dashboard. Click **Assignments** at the top of the screen to view user progress. Click **Dataflows** to view dataflow progress. Use the scale to choose the metric for the scale (day, week, month) and the number of units to display.

Assigning Records

The **Maintenance > Assignment** section under on the **All Records** page enables you to reassign records from one user to another.

1. On the Data Stewardship Portal page, click **Manage**.
This displays the **Manage** page.
2. Click **Assignments** to expand the assignment options.
3. In the **User** field, select a user whose exceptions you want to assign to another user.
4. To reassign all records belonging to a user, skip to step 6 on page 197. To reassign a portion of a user's records, complete one or more of these fields:

Field	Description
Dataflow name	The name of the dataflow producing the records.
Stage label	The name of the stage producing the records.
Job ID	The ID assigned to the job containing the records.

Field	Description
Data domain	The kind of data assigned in the Exception Monitor.
Quality metrics	The kind of metric assigned in the Exception Monitor.
From date/time	The start date in a range of dates in which the records were created.
To date/time	The end date in a range of dates in which the records were created.
Record status	Select one or more status from the list. If none is selected, all records matching other filters are reassigned. For example, select status of Edited and Rejected records and only records with those status values are reassigned.

5. After completing your selections in the **User**, **Dataflow name**, and **Stage label** fields (at minimum), you can further refine the filter based on exception field values.
- Click the Add expression button **+**.
 - In the **Field Name** column, select the field you want to filter on.
 - In the **Operator** column, select one of the following:

Operator	Description
Is equal to	Looks for records that have exactly the value you specify. This can be a numeric value or a text value. For example, you can search for records with a MatchScore value of exactly 82, or records with a LastName value of "Smith."
Is not equal to	Looks for records that have any value other than the one you specify. This can be a numeric value or a text value. For example, you can search for records with any MatchScore value except 100, or records with any LastName except "Smith."
Is greater than	Looks for records that have a numeric value that is greater than the value you specify.
Is greater than or equal to	Looks for records that have a numeric value that is greater than or equal to the value you specify. For example, if you specify 50, you would see records with a value of 50 or greater in the selected field.

Operator	Description
Is less than	Looks for records that have a numeric value that is less than the value you specify.
Is less than or equal to	Looks for records that have a numeric value that is less than or equal to the value you specify. For example, if you specify 50, you would see records with a value of 50 or less in the selected field.
Contains	Looks for records that contain the value you specify in any position within the selected field. For example, if you filter for "South" in the AddressLine1 field, you would see records with "12 South Ave.", "9889 Southport St.", "600 South Shore Dr.", and "4089 5th St. South".
Does not contain	Looks for records that do not contain the value you specify in any position within the selected field. For example, if you filter for "South" in the AddressLine1 field but select "does not contain," you would not see records with "12 South Ave.", "9889 Southport St.", "600 South Shore Dr.", and "4089 5th St. South".
Starts with	Looks for records that start with a particular value in the selected field. For example, if you filter for "Van" in the LastName field you would see records with "Van Buren", "Vandenburg", or "Van Dyck".
Ends with	Looks for records that end with a particular value in the selected field. For example, if you filter for records that end with "burg" in the City field, you would see records with "Gettysburg", "Fredricksburg", and "Blacksburg".

- d) In the **Value** column, enter the value to use as the filtering criteria.

Note: The search value is case-sensitive. This means that searching for SMITH will return only records with "SMITH" in all upper case, but not "smith" or "Smith."

- e) To filter on more than one field, add multiple filters by clicking the add field filter icon again. For example, if you want all records with a LastName value of "SMITH" and a State value of "NY" you could use two filters, one for the LastName field and one for the State field.

This example would return all records with a value of "FL" in the State field:

The screenshot shows a filter configuration interface. At the top left, there are icons for adding (+), deleting (trash), and refreshing (refresh) filters. Below these icons is a table with three columns: 'Field name', 'Operator', and 'Value'. The first row of the table has a dropdown menu for 'Field name' with 'State' selected, a dropdown menu for 'Operator' with 'is equal to' selected, and a text input field for 'Value' containing 'FL'.

This example would return all records that do not have a PostalCode value of 60510:

Field name	Operator	Value
StateProvince	Is equal to	NY
PostalCode	Is not equal to	14226

This example would return all records with a StateProvince of "NY" with all postal codes except 14226.

Field name	Operator	Value
PostalCode	Is not equal to	60510

6. Click **Reassign**.
7. Optional: Select the number of exceptions you want to reassign.
You can assign all or some of the resulting exceptions to the new user. For example, if you enter "10" as the limit, only the first 10 records meeting the criteria will be reassigned and the remaining records meeting the criteria will not be reassigned.
8. Select another user in the **Reassign** box.
9. Click **Confirm**.

Deleting Records

Occasionally you may want to delete records from the exception repository. For instance, you could have residual records from testing the system or records that were mistakenly considered exceptions after processing, or you may want to process and delete approved records first and then re-run the same job again. The **Purge** section, under **Maintenance**, on the **All Records** page enables you to do this.

Note: You can also use the `bsm delete exceptions` command in the Command Line Interface (CLI) to delete records. For details about the CLI commands, see the Administration Utility section of the *Spectrum Technology Platform Administration Guide*.

1. In the Data Stewardship Portal, click **Manage**.
This displays the **Manage** page.
2. Under **Maintenance**, expand the **Purge** section.
3. In the **Dataflow name** box, select a dataflow.
4. In the **Job ID** box, either select a job ID or **All**.
Select **All** to remove records from every job in the selected workflow.
5. Optional: To remove all performance data, you can check the **Remove data quality report data** check box.

If you do not check this check box, the records will be removed from the repository but performance data will still appear on the Data Quality page. This option is unavailable if you choose to add filter criteria in step 6 on page 198.

6. Optional: You can add filter criteria to specify records to be deleted.
 - a) In the **Stage label** box, select the stage that produced the records.
 - b) In the **Approval status** box, select one or more status from the list.
If none are selected, all records matching other filters are deleted.
 - c) For each filter you want to apply, click the Add expression button **+**.
In each case, specify a **Field name**, an **Operator**, and **Value** to filter based on a field values. You must have specified a **Job ID** and **Stage Label** in order to define field expressions in this step.
7. Click **Remove**.

The Data Quality Page

The Data Stewardship Portal **Data Quality** page provides information on trends among records.

Note: This page will be deactivated if you turn off data quality reporting on the Data Stewardship Settings page of Management Console.

Identifying Trends

The **Data Quality Trends** page depicts the following statistical information by dataflow and domain:

- Total number of conditions processed
- Total number of exception conditions
- Percentage of conditions that were processed successfully
- The trend of your data in 30-day intervals

This information can be further broken down by stage if you select a dataflow or by metric if you select a domain. The values that appear here are determined by dataflows with Exception Monitor stages that have been run over the last 30 days as well as their respective domains and metrics.

- If you are viewing information by domain, select a **Dataflow name** if you want to view information for a specific dataflow and a **Stage label** if you want to view information for a specific stage.

Note: You must select a single dataflow if you want to also filter the results based on a stage. Otherwise, you will see data for all dataflows.

- If you are viewing information by dataflow, select a **Domain** if you want to view information for a specific domain and a **Metric** if you want to view information for a specific metric. Otherwise, you will see data for all domains.
- Select a duration for the **Scale** to specify how far back you want the data to go. You can select any number of days, weeks, or months.
- Select a dataflow if you want to break out results by stage.
- Select a domain if you want to break out results by data quality metric.

As you make selections to view different parts of the data, you will notice that charts appear reflecting the status of that particular data. For instance, if you look at domain data, you could narrow that data down by a specific domain (such as "Address") and then by a particular metric (such as "Completeness"). The charts will update accordingly as you make different selections.

How to Video - Approval Flow in Data Stewardship

This video explains how you can build Approval Flow in Data Stewardship.

<https://www.youtube.com/embed/ymQTa35L1QQ>

Introduction

The Data Stewardship Portal provides tools for reviewing, modifying, and approving records that failed based on conditions in the Exception Monitor stage. Use the Data Stewardship Portal to manually enter correct or additional data in a record. Typical exceptions are:

- Address verification failures
- Geocoding failures
- Low-confidence matches
- Merge/consolidation decisions

For example, if a customer record fails an address validation process, you could use the search tools to conduct research and determine the customer's address, then modify the record so that it contains the correct address. The modified record could then be approved and reprocessed by Spectrum Technology Platform, sent to another data validation or enrichment process, or written to a database, depending on your configuration. You could also use the Data Stewardship Portal to add information that was not in the original record.

Data Stewardship Portal consists of the following components.

- **Dashboard**—The Dashboard displays the number of records assigned to the user, counts of resolved records that are still assigned to the user, counts of records that are still unresolved, and the daily performance (number of records accepted/rejected). Counts are broken down by Dataflow and Label.
- **Editor**—The editor allows you to modify records and change status of records within review cycles.
- **Manage**—A browser-based tool that enables you to review and manage record activity for all users, including reassigning records from one user to another. Summary and detail charts show the number of records with resolved and unresolved status, and the number of records assigned to each user.
- **Data Quality**—A browser-based tool that provides information on trends for records and enables you to identify key performance indicators and send notifications when certain conditions are met.

Accessing the Data Stewardship Portal

Before you perform this procedure, you must have a user name and password required to access the Spectrum server on your network. Login credentials are normally provided by the Spectrum Technology Platform administrator.

You can open Spectrum and the Data Stewardship Portal using either the Windows Start menu or a web browser.

- [Use the Windows Start menu to access the Data Stewardship Portal](#) on page 200
- [Access the Data Stewardship Portal from the Spectrum Home page](#) on page 201
- [Access the Data Stewardship Portal directly](#) on page 201

All of these methods open the Data Stewardship Portal dashboard.

Use the Windows Start menu to access the Data Stewardship Portal

You can use this Windows Start menu if you are working on the computer that is running the Spectrum Technology Platform server.

1. Click **Start > Precisely > Welcome Page**.
2. Click **Spectrum Data Quality > Data Stewardship Portal > Open the Data Stewardship Portal**.
3. On the **Sign in** page, enter your Spectrum user name and password to sign in.

Access the Data Stewardship Portal from the Spectrum Home page

To perform this procedure, you need to know the name of the host server running Spectrum Technology Platform and the port number used by Spectrum Technology Platform.

1. Open a web browser and navigate to `http://SpectrumServer:port/`.
For example,
`http://myserver:8080`
2. Click **Spectrum Data Quality > Data Stewardship Portal > Open the Data Stewardship Portal**.
3. On the **Sign in** page, enter a valid user name and password to sign in.

Access the Data Stewardship Portal directly

To perform this procedure, you need to know the name of the host server running Spectrum Technology Platform and the port number used by Spectrum Technology Platform.

1. Open a web browser and navigate to `http://SpectrumServer:port/bsmportal`.
For example,
`http://myserver:8080/bsmportal`
2. On the **Sign in** page, enter a valid user name and password to sign in.

The Dashboard page

The Dashboard page in Data Stewardship Portal displays data that summarizes the status of records belonging to you and other users.

The Dashboard page is displayed by clicking **Dashboard** on the **Data Stewardship Portal** page. The Dashboard displays status and statistics for a selected user. The default is the currently logged in user. You choose to view data for other users if you have modify permissions. A different user may be selected in the box at the top of the page. Only users who currently have exceptions assigned to them appear in users box.

Note: The progress charts for user and dataflow will not appear if you have turned off progress tracking on the Data Stewardship Settings page of Management Console.

Total	Shows total number of records for the currently selected user.
Remaining	Shows number and percentage of records that remaining for the currently selected user.
Resolved	Shows the number and percentage of records that have been resolved for the currently selected user.
Status	The chart shows the selected user's percentages of resolved and remaining records.
Daily Progress	<p>The selected user's progress by day, week, or month. This shows number of records accepted by the user. For an approval flow, it shows the number of records accepted and rejected by the user.</p> <ul style="list-style-type: none"> • Click the Display table button  to view progress in table format. • Click the Display chart button  to view progress in chart format. • Click the Options button  to select the duration of Day, Week, or Month.
Breakdown by dataflow	This section lists active dataflows. To edit a dataflow on the Data Stewardship Portal Editor page, click the dataflow name. To filter the list, enter a portion or all of the characters in a dataflow name in the Filter box  .
Resolved Timeline	This section shows the resolved status of approval flows. Information includes, the ratio and percentage of resolved records, the number of records waiting to be resolved or accepted, the number of records resolved, and the total number of records.

The Data Stewardship Portal Editor

The Data Stewardship Portal **Editor** enables users to manually review, modify, and resolve records that cannot be corrected by stages in a dataflow. Users can also use the Data Stewardship Portal **Editor** to resolve duplicate records and use search tools to look up information that assists in editing, resolving, and rerunning records. A user with modify permission is able to correct incorrect or missing data in a record and resolve a record. Resolved records can then be read from the Data Stewardship repository using the Read Exceptions stage for additional processing.

When an approval flow is defined for data, the Exception Monitor assigns records to users in a sequence of one or more acceptance levels. Exception Monitor initially assigns records to the default user, referred to as the *data steward*. The data steward reviews and edits a record before accepting it. After the data steward accepts a record, it is assigned to the next level of the approval flow. In each level after the data steward, a user can modify, accept, or reject a record. A record that is rejected at any level is assigned back to the data steward. After the last-level user approves and

saves a record, it is set to resolved. The resolved records can then be read by the Read Exceptions stage for additional processing and are purged from the Data Stewardship repository.

Customizing Data Stewardship Portal Editor Contents

There are multiple ways you can customize what is shown on the Data Stewardship Portal **Editor** page. You select settings on the side panel to view records in the exception repository for a particular **User** (if the current user has access to other user exception assignments), **Dataflow name**, **Stage label**, **Job ID**, **Status**, **Date**, or **Condition** or **Data Domain** or **Metrics**. You can use the Set filter  tool to only display records that meet certain criteria. The Configure view  options allow you to show or hide fields in the editor.

Do not forget that after you configure settings described in the previous paragraph in combination with the **items per page** setting, there may still be multiple pages of results that are not shown on the initial screen. For example, you set **items per page** to 10 and then apply a field filter to return only the records that have a specific postal code. You may initially see only 10 records, but should you check the bottom right corner, as it may show that there are additional items in the repository.

Search Options

Search options determine which records are visible in the Exception **Editor**.

Click the **SEARCH** tab on the left side of the Data Stewardship Portal **Editor** page to expand or collapse search options. Click any of the following headings on the **SEARCH** tab to expand or collapse criteria for an option.

User	The ID of the user assigned to the dataflow. This information will be visible only if you have modify permissions. This selection is required.
Dataflow name	The name of the dataflow that generated the records. This selection is required.
Stage label	The user-defined name given to the Exception Monitor stage in the dataflow. This information is particularly useful in cases where a dataflow contains multiple Exception Monitor stages. If the person who created the dataflow gave each Exception Monitor stage a meaningful label you can identify which Exception Monitor produced the record. The default label is "Exception Monitor". This selection is required.
Job ID	A numeric identifier assigned to a job by the system. Each time a job runs it is assigned a new job ID. You can choose to view records for all jobs or select a specific job.
Status	Lists possible approval status categories for records in the exception repository. The number of records currently in a category is displayed next to the status label. These values are updated each time a user saves changes. Check the check box next to a status label to display records that currently display the status in the Record status field. You can check multiple check boxes on this panel.

- Date** Choose to display records from dataflows run on specified dates. You can specify either or both start and end dates. Click the date button  select a date. Click the time button  to specify the time. The default is 12:00 AM.
- Condition**
- **Name**—Choose a condition name to only view records that were created as a result of a particular condition. If you specify a name, then both **Data Domain** and **Quality Metrics** are set to **All**.

Note: To display condition names for each record, see [Search Options—Condition Name](#).
 - **Data Domain**—Specifies the kind of data to which the record belongs, such as Customer or Product. Examples of data domains include Name, Address, and Phone Number. This selection helps you identify which fields in the record require editing. If you specify a data domain, then **Name** is set to **All**.
 - **Quality Metrics**—The quality measurement that failed the record. Examples of quality metrics include Accuracy, Completeness, and Uniqueness. This information helps you determine why the record was identified as an exception. If you specify a quality metric, then **Name** is set to **All**.

Field filters

Field filters control limit records that are displayed in the Data Stewardship Portal **Editor** according to field values.

After records are loaded into Data Stewardship, you filter fields to display only those records in which you are interested. The Data Stewardship Portal initially displays all records from a single Spectrum Technology Platform dataflow. Field filters display only those records that meet specified criteria for one or more fields. Once filtering is defined for a dataflow, it is applied each time you open the dataflow on the **Editor** page until you clear the filters.

Note: You must apply any selection options described in [Search Options](#) on page 167 before using the field filter tool.

Filters are set by clicking the Set filter button  on the toolbar. Doing so displays the Filter icon  on field headings for fields that can be filtered. You can click the Filter icon  on a field heading to set a filter for the field. After you set a filter on a field, its Filter icon  displays a blue background:

City	FirstName
KEENE	LAKSHMI
KEENE	LAKSHMI

Filters determine which records are visible in both Tabular view and Form view, although they can only be created in the Tabular view. The filter icon on the Form view status bar signifies that records are currently filtered:



Related concepts

[Lookup Tables \(Data Quality Guide\)](#) on page 150

Related tasks

[Setting field filters](#) on page 169

This task describes how to filter the list of records according to field values.

Related reference

[Field filter operators](#) on page 170

This topic describes operators that you may apply in a field filter. The available operators depend on whether a field is a text, date, or numeric field.

Setting field filters

This task describes how to filter the list of records according to field values.

1. On the Data Stewardship Portal **Editor** page, make sure that you are in Tabular view .
2. Click the Set filter button  on the toolbar.
This displays the Filter icon  on every field heading that can be filtered.
3. Click the Filter icon  on a field heading.
4. In the **Show items with value that** list, select the operator that you want to use.
For more information, see [Field filter operators](#) on page 170.
5. Enter one or more values for the filter.
If a lookup table was defined in the Data Stewardship settings and is associated with the column, click the value box to expand the list and select values. Otherwise, you can type and add values to the value box.
6. Click the **Filter** button to apply the criteria.
Only records whose data matches the filter criteria are then visible. If you add multiple values to a filter, records are displayed whose data matches the criteria for any of the values. All records are hidden if there is no match.
7. To filter on additional fields, repeat steps [3](#) on page 205 through [6](#) on page 205 for each field that you want to filter.
When you set filters on multiple fields, records are displayed whose data matches the criteria set for all of the fields.

You can click the Set filter button  again to clear filters currently set on fields.

Related concepts

[Lookup Tables \(Data Quality Guide\)](#) on page 150

Related tasks

[Clearing field filters](#) on page 170

You can clear a single field filter or all field filters.

Related reference

[Field filter operators](#) on page 170

This topic describes operators that you may apply in a field filter. The available operators depend on whether a field is a text, date, or numeric field.

Clearing field filters

You can clear a single field filter or all field filters.

- To clear the filter from a single field in Tabular view, first click the Filter icon  on the field heading, then click the **Clear** button.
- To clear filters from all fields in either Tabular view or Form view, click the Set filter button  on the toolbar.

Related tasks

[Setting field filters](#) on page 169

This task describes how to filter the list of records according to field values.

Field filter operators

This topic describes operators that you may apply in a field filter. The available operators depend on whether a field is a text, date, or numeric field.

Is equal to	Looks for records that have exactly the value you specify. This can be a numeric value or a text value. For example, you can search for records with a MatchScore value of exactly 82, or records with a LastName value of "Smith."
Is not equal to	Looks for records that have any value other than the one you specify. This can be a numeric value or a text value. For example, you can search for records with any MatchScore value except 100, or records with any LastName except "Smith."
Starts with	Looks for records that start with a particular value in the selected field. For example, if you filter for "Van" in the LastName field you would see records with "Van Buren", "Vandenburg", or "Van Dyck".
Contains	Looks for records that contain the value you specify in any position within the selected field. For example, if you filter for "South" in the AddressLine1 field, you would see records with "12 South Ave.", "9889 Southport St.", "600 South Shore Dr.", and "4089 5th St. South".
Does not contain	Looks for records that do not contain the value you specify in any position within the selected field. For example, if you filter for "South" in the AddressLine1 field, you would not see records with "12 South Ave.", "9889 Southport St.", "600 South Shore Dr.", and "4089 5th St. South".
Ends with	Looks for records that end with a particular value in the selected field. For example, if you filter for records that end with "burg" in the City field, you would see records with "Gettysburg", "Fredricksburg", and "Blacksburg".

Is greater than	Looks for records that have a numeric value that is greater than the value you specify.
Is greater than or equal to	Looks for records that have a numeric value that is greater than or equal to the value you specify. For example, if you specify 50, you would see records with a value of 50 or greater in the selected field.
Is less than	Looks for records that have a numeric value that is less than the value you specify.
Is less than or equal to	Looks for records that have a numeric value that is less than or equal to the value you specify. For example, if you specify 50, you would see records with a value of 50 or less in the selected field.
is after	Looks for records that have a date or time value that is later than the value you specify.
Is after or equal to than	Looks for records that have a date or time value that is equal to or later than the value you specify.
Is before	Looks for records that have a date or time value that is earlier than the value you specify.
Is before or equal to than	Looks for records that have a date or time value that is equal to or earlier than the value you specify.

Related tasks

[Setting field filters](#) on page 169

This task describes how to filter the list of records according to field values.

Viewing Records

To open the Data Stewardship Portal **Editor** to view records, click **Editor** on the Data Stewardship Portal menu bar. The Data Stewardship Portal **Editor** will also open when you click a dataflow in the **Breakdown by dataflow** list on the Dashboard.

On the **Editor** page, you can view records in either Tabular view or Form view.

Tabular view

In Tabular view you can load up to 100 records per page. You can change how many records are displayed per page in the **items per page** box at the bottom of the screen. In this view, you can scroll through the list and edit records in any order. When you finish editing records on a page, you can click the Save changes button  to simultaneously save changes to all of the records on a page.

Form view

In Form view you view and edit data for one record at a time. You can click the navigation buttons to scroll between records in the current dataflow. In this view you must save the edits to a record

before you proceed to a different record. You cannot save edits more than one record at a time in Form view.

Viewing Record Details

Regardless whether you use Tabular view or Form view, the Exception grid shows all the fields for a record as well as record status, exception type, and any comments that have been added to the record. In the Tabular view, you can view additional details about a record by clicking the arrow in the left column of a record. In the Form View, click **Show detail** beneath the form fields. In either view, expanding details for a record shows the following tabs.

Details tab on page 172	This tab shows details about a record.
Condition tab on page 173	The tab shows information about the condition that identified the record as an exception.
History tab on page 173	This tab shows historical information for a record.
Approval flow tab on page 174	This tab shows the current status of a record and the levels that have been completed in an approval flow.

Details tab

This tab shows details about a record.

Job ID	A numeric identifier assigned to a job by the system. Each time a job runs it is assigned a new job ID.
Group by	If the dataflow was configured to return all records in the records group, this shows the field by which the records are grouped. This only applies to dataflows that perform matching, such as dataflows that identify duplicate records or dataflows that group records into households.
Created on	The date and time when the Exception Monitor identified the record as an exception.
Record type	The type of record. This is one of the following: <ul style="list-style-type: none"> • Exception • Group exception • Group non-exception • Best of Breed
Resolved	Shows whether record is resolved (Yes or No).
Approval flow types	Shows inclusion in approval flow types.

Condition tab

The tab shows information about the condition that identified the record as an exception.

Name	The name of the condition that identified the record as an exception. Condition names are defined by the person who set up the dataflow.
-------------	--

Domain	The kind of data that resulted in an exception. Examples of data domains include Name, Address, and Phone Number. This information helps you identify which fields in the record require editing.
Metrics	The quality measurement that the record failed. Examples of quality metrics include Accuracy, Completeness, and Uniqueness. This information helps you determine why the record was identified as an exception.

History tab

This tab shows historical information for a record.

Every time a record is changed, field values and historical information are retained by the system. The table on this tab shows saved versions of a record. The column heading for a field is highlighted in blue when its value was changed at any time in the history of the record. In these highlighted columns, a field value is highlighted in blue when it does not match the value shown in the preceding version of the record.

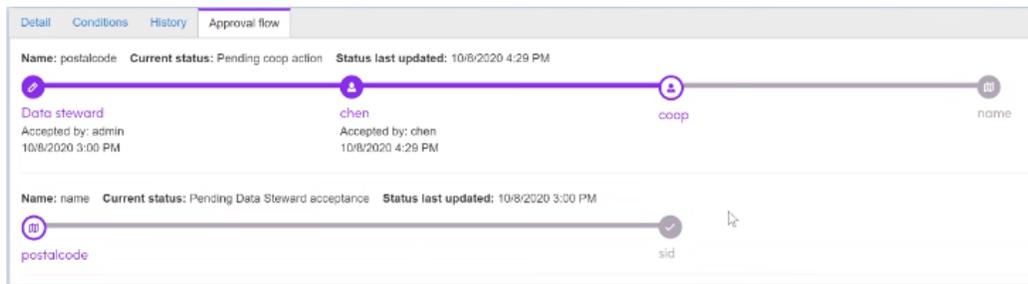
The initial columns in the table provide the following information about each version of a record:

Version	The revision number of the change. This starts with 1.0 for the original version of a record and increments 1.1, 1.2, 1.3, and so forth each time that changes are saved to a record.
Modified time	Date and time when record was last saved.
Modified by	The user who made the last change.
Assigned to	The user to whom the record was assigned at the time of its revision.
Record status	The current status of the record. This is one of the following: <ul style="list-style-type: none"> • Accepted—Record has been accepted in an approval flow. • Edited—Record has been edited. • Failed—Record failed revalidation. • New—Record has not been reviewed or edited. • Pending—Record has been reviewed or edited but has been retained for additional review. • Rejected—Record has been rejected in an approval flow. • Resolved—Accepted record has been saved for additional processing.
Condition name	The condition in the Exception Monitor stage that generated the record.
Comment	The most recent comments that were entered for the record.

Approval flow tab

This tab shows the current status of a record and the levels that have been completed in an approval flow.

Shows the following information for an approval flow. When the record is in multiple approval flows, the information and visual representation is provided for each flow. When a record is flagged by more than one condition, it progresses through each approval flow. Each status line represents a single approval flow plan.



Name	The name of the condition.
Current status	Shows the whether a record is pending action by a reviewer, rejected, accepted, or resolved.
Status last updated	Shows the date and time that the status was last updated.

Sorting records

In Tabular view, you can sort records by values in a single column. You can sort records by text (A to Z or Z to A), numbers (smallest to largest or largest to smallest), and dates and times (oldest to newest and newest to oldest).

1. In Tabular view, choose the field that you want to use to sort records.
2. Click the column heading to sort records based on ascending order of values in the column. This displays the ascending sort  icon next to the column heading to indicate that records are sorted based on ascending values in the column.
3. To sort based on a descending sort of values in the column, click the column heading again. This displays the descending sort  icon next to the column heading to indicate that records are sorted based on descending values in the column.

You can sort on any column except the **Condition name** column, which is displayed in the record table when the **Condition name** field check box is checked in the Configure view. For more information, see [Configuring the view](#) on page 175.

Grouping records

In Tabular view you can group records by values in one or more columns. When you group by values in more than one column they are grouped hierarchically in the order that you choose.

To complete this task, the **Allow grouping records list** check box must be checked in the **Configure view** options. For more information, see [Configuring the view](#) on page 175.

1. For the first column on which to group records, drag the column heading to the panel above the column headings.

Record status	Type	Comment	Condition name	AddressLine1	City	CollectionNum
City: ABINGTON						
> New	🚩		MatchScore	2516 PEERSHING ...	ABINGTON	0
City: AMBLER						
> New	🚩		MatchScore	202 SPOUT ROAD	AMBLER	0

Completing this step groups records by values within this column.

- Optional: To change the sort order of values in the group, click the ascending ▲ or descending ▼ icon on the column entry in the panel.

This changes the order for records displayed on a page.

Note: To group records and maintain sort order across pages, sort the column before adding it to the grouping. For more information, see [Sorting records](#) on page 174.

- To add subordinate groupings, you can drag additional column headings to the panel above the table.

Completing this step creates groupings in the order that you add headings.

- Optional: To change the order in which records are grouped, drag a column heading already in the panel before or after another column heading in the panel.

Configuring the view

You can select which fields appear and change the order in which they appear by clicking the Configure view button ⚙️. This displays the **Configure view** options.

Changes to the **Configure view** settings are saved for a user and dataflow name, so that when you reopen a dataflow the configuration from the previous session is applied. Changes made here apply to editing records in both Tabular view and Form view.

Filter Enter all or part of a field name in the Filter box 🔍 to only display field names that contain the entered characters in their names. The search is case insensitive and the text can be anywhere in the name of a field.

All List all fields. Check or uncheck a field to expose or hide it in both Tabular view and Form view.

Visible Lists fields that are currently checked in the **Field name** column.

Hidden Lists fields that are currently not checked in the **Field name** column.

Display Condition name field In the **Field name** column, check a check box to display a column in both Tabular view and Form view. This field displays the name of the condition that caused the record to be added to the repository. Records that meet criteria for multiple conditions will display condition names as a comma separated list.

Note: To only display records for a single condition, see the [Search Options—Condition Name](#).

Change field order The fields are listed in the **Configure view** dialog box in the same order that they are displayed in the records table. You can change the order in which fields appear by dragging and dropping fields in the desired order. However, you cannot rearrange fields while there is text in the Filter box . You must first clear the contents of the Filter box before you can reorder fields.

Note: You can also order fields in Tabular view. Drag column headings on the record table to reorder the columns. Any changes to the table itself will also appear in the **Configure view** dialog box the next time that you open it.

Configure field widths You can specify the width of columns that appear in the Form view by entering values in the **Form columns** boxes. Values for adjacent visible fields in this column must sum to 12. For example, if you wanted four fields to be placed side-by-side in four equal-length fields on a single row, you could give each field a column width of "3." Likewise, if you wanted to include just two fields of equal length in a row, you could enter "6" for both fields under **Form columns**. These proportional settings are illustrated in the first figure, [Figure 3: Form columns configured for a rowupdate image, checkboxes changed, first option changed to "Optimize record editors for keyboard navigation"](#) on page 213. The fields in the figure are set to 4, 3, 2, and 3 to define four adjacent fields of various widths that sum to 12. The second figure, [Figure 4: Fields in Form view rendered as configured](#) on page 213, shows the result from Form view.

Configure view ×

▼

All
Visible
Hidden

Field name	Form columns
↕ <input checked="" type="checkbox"/> Approved	Default ▼
↕ <input checked="" type="checkbox"/> Type	Default ▼
↕ <input checked="" type="checkbox"/> Comment	Default ▼
↕ <input checked="" type="checkbox"/> AddressLine1	4 ▼
↕ <input checked="" type="checkbox"/> City	3 ▼
↕ <input checked="" type="checkbox"/> State	2 ▼
↕ <input checked="" type="checkbox"/> PostalCode	3 ▼

Optimize tabular editor for keyboard navigation
 Allow grouping records list

OK
Cancel
Reset to default

Figure 3: Form columns configured for a row

Approved
 Comment

555.55HC71 BOX60

CTR STRAFFORD

NH

Figure 4: Fields in Form view rendered as configured

Note: When you configure field widths under **Format columns**, the configured fields will no longer be responsive to the size of the browser window. In contrast, fields with the **Default** setting are always responsive to the size of browser window, wrapping and unwrapping when the window is resized.

Enable Keyboard Navigation

Check the **Optimize tabular editor for keyboard navigation** check box to use the Tab key to navigate records in the Tabular view of the Data Stewardship Portal **Editor** page.

Group records in Tabular view

Check the **Allow grouping records list** check box to group records by their column values. This displays a panel above the column headings in which to configure groupings. Records may be grouped hierarchically in descending or ascending order on a page by values in one or more columns. This check box is checked by default. For more information, see [Grouping records](#) on page 175.

Reset to Default

Click the **Reset to default** button to check all of the check boxes in the **Field name** column and to set **Form columns** to "Default".

Editing and resolving records

The purpose of editing a record is to correct or add information and then accept and resolve the record so that it can be processed successfully. Editing a record may involve using other Spectrum Technology Platform services or consulting external resources such as maps, the Internet, or other information systems in your company.

After reviewing records, you can edit, accept, and resolve them directly on the Data Stewardship Portal **Editor** page. You can edit one or more records simultaneously in the Tabular view, and you can edit one record at a time in the Form View.

Note that read-only fields cannot be edited. If you want to make a read-only field editable, you would need to delete all records for that dataflow and job ID and run the dataflow again after configuring the fields accordingly in the Write Exceptions stage. This would produce new records with editable fields. Furthermore, you cannot edit a field to contain a value that does not match its data type. For example, you cannot enter alphabetic characters in a field that has a numeric data type.

Editing records in Tabular view

Perform steps in this task to edit and approve records in Tabular view.

1. In Data Stewardship Portal, open records for the dataflow that you want to edit on the **Editor** page.
To do this, click **Editor** on the Data Stewardship Portal home page, then in the **Dataflow name** box on the sidebar, select the name of the dataflow.
2. Make sure that the records are displayed in Tabular view.
Click the Tabular view button  on the **Editor** toolbar.
3. You can now edit field values for a single record or for multiple records:
 - [To edit a field for a single record in Tabular view](#) on page 214
 - [To edit a field for multiple records in Tabular view](#) on page 215

To edit a field for a single record in Tabular view

1. In any record, click to edit a value.
Right-click the field to access **Cut**, **Copy**, and **Paste** commands. A field will be grayed out if it cannot be edited.

To undo a change, select the record you want to undo and click the Undo changes button  on the toolbar.

After you edit a field, green triangle  in the upper-left corner that indicates edits are not yet saved. After you save changes to the record, the status changes to Edited.

2. Optional: You can add comments about your changes in the **Comment** box.
Comments are visible to other users and can be used to help keep track of the changes made to the record or to recommend additional changes.
3. After you are satisfied with corrections made to a record, click Accept .
This changes the **Record status** to Accept.

Note: Anytime before you click the Save button, you can click Accept  or Reject  to change the status from Accept to Pending. This allows you to continue editing the record at a later time.

4. Click the Save changes button  when you are finished editing records.
This changes status for the record to Resolved. The records may not appear in the table if they are no longer available or no longer match the search/filter criteria.

To edit a field for multiple records in Tabular view

1. Click the field in the first record that you want to change.
2. Press Shift to extend the selection to multiple records or Ctrl to add a single record to the selection, and click the same field in another record.
For instance, if you want to change all instances of "L.A." to say "Los Angeles," click within the City field for one of the records. Then press Ctrl or Shift and click within that same column to extend or add to the selection.
3. Edit the field value in the last record that you selected.
Right-click the field to access **Cut**, **Copy**, and **Paste** commands. A field will be grayed out if it cannot be edited.

To undo a change, select the records you want to undo and click the Undo changes button  on the toolbar.

After you edit a field, green triangle  in the upper-left corner that indicates edits are not yet saved. After you save changes to the record, the status changes to Edited.

4. Optional: You can add comments about your changes in the **Comment** box.
Comments are visible to other users and can be used to help keep track of the changes made to the record or to recommend additional changes.
5. After you are satisfied with corrections made to a record, click Accept .
This changes the **Record status** to Accept.

Note: Anytime before you click the Save button, you can click Accept  or Reject  to change the status from Accept to Pending. This allows you to continue editing the record at a later time.

- Click the Save changes button  when you finish editing records on a page. This changes status for the records to Resolved. The records may not appear in the table if they are no longer available or no longer match the search/filter criteria.

Editing fields in Form view

Perform steps in this task to edit records in the Form view.

- In the Data Stewardship Portal **Editor** page, display records in the exception repository for the dataflow that you want to edit .
To do this, click **Editor** on the Data Stewardship Portal toolbar, then in the **Dataflow name** box on the sidebar select the name of the dataflow.
- Make sure that the records are displayed in Form view.
Click the Form view button  to display fields in the first record.
- In any field, click to edit a value.
Right-click a field to access **Cut**, **Copy**, and **Paste** commands. A field will be grayed out if it cannot be edited.

To undo a change, select the record you want to undo and click the Undo changes button  on the toolbar.
An edited field is highlighted green to indicate unsaved changes. After you save changes to the record, the status changes to Edited.
- Optional: You can add comments about your changes in the **Comment** box.
Comments are visible to other users and can be used to help keep track of the changes made to the record or to recommend additional changes.
- After you are satisfied with corrections made to a record, click Accept .
This changes the **Record status** to Accept.

Note: Anytime before you click the Save button, you can click Accept  or Reject  to change the status from Accept to Pending. This allows you to continue editing the record at a later time.
- Click the Save button .
This changes status for the records to Resolved. You will either see the same record containing your changes or you will see the next record in the list because the record you changed is no longer available or no longer matches the search/filter criteria.
- Use the navigation buttons at the bottom of the form to go to previous or next record.
You can also go directly to the first or last record.

Exporting Records

You can export data for all records from a dataflow to a comma separated value (CSV) file. You can then import the data into other applications such as Microsoft Excel.

This procedure saves data for all of the records from a dataflow to a CSV file. The first (or header) row in the CSV file lists the field names. Subsequent rows list the corresponding values for each record. Field names and data values are enclosed between double quotation marks (").

1. In Data Stewardship Portal, open records for the dataflow that you want to export on the Data Stewardship Portal **Editor** page.

To do this, click **Editor** on the Data Stewardship Portal toolbar, then in the **Dataflow name** box in the sidebar select the name of the dataflow.

2. Make sure that the records are displayed in Tabular view.

If necessary, click the Tabular view button  on the **Editor** toolbar.

3. On the **Editor** toolbar, click the Export all records button .

4. If prompted, choose to save the file.

All the records created by the dataflow are exported to the browser default download folder. The exported file name is *DataflowName_StageLabel.csv*, with the dataflow name and stage label of the dataflow that created the records. If records created by a dataflow are exported repeatedly, the files are numbered incrementally in parentheses after the dataflow name as

DataflowName_StageLabel(number).csv.

After you complete this procedure, you can open the file and edit records in another application, such as Microsoft Excel or a text editor. After you complete the changes, you can load the data back into a workflow in Enterprise Designer or Flow Designer.

Approval flows

Approval flows assign records to sequential levels for review and acceptance by a user.

Records are initially assigned to the data steward specified in the Exception Monitor stage. The data steward can review, edit, and accept records. If a record is accepted, it is then assigned to the first approver as defined in the approval flow plan. At this level, the approver can Accept or Reject a record. If a record is accepted it is be assigned to the next approver in the approval flow. If a record is rejected, it is assigned back to the data steward. This is repeated at subsequent levels until the final approver accepts and resolves the record. Resolved records can then be read using the Read Exception stage for additional processing.

Creating approval flows

This process defines tasks that must be completed to create an approval flow. The approval flow process is illustrated in the [Approval Flow Process diagram](#).

- [Define approval flow types](#) on page 218

Define approval flow types in the Data Stewardship Settings of the Management Console.

- **Specify the approval flow type for a condition** on page 218

Assign an approval flow type to a condition to be used in an approval flow. This task is performed in the Exception Monitor stage. Perform this procedure in the Exception Monitor settings in an Enterprise Designer dataflow.

- **Enable approval flow types in security roles** on page 219

Execute permission must be applied to an approval flow type within a security role before you can enable it in an approval flow. Perform this procedure in the Roles settings of the Management Console.

- **Assign approval flow role to new or existing users** on page 219

To define reviewers in an approval flow, assign an approval flow role to new or existing users. Perform this procedure in the Users settings of the Management Console.

- **Define an approval flow** on page 220

Define levels and specify the role and user for each level in the Management Console Approval Flow settings. Perform this procedure in the Approval Flows settings of the Management Console.

Define approval flow types

Define approval flow types in the Data Stewardship Settings of the Management Console.

1. Open the Spectrum Technology Platform **Management Console**.
On the Spectrum Technology Platform **Home** page, click **Platform Client tools > Web > Open Management Console**.
2. On the **Management Console** page, click **Resources > Data Stewardship Settings**.
3. On the **Data Stewardship Settings** page, click **Approval Flow Types**.
4. Click the Add type button  to create a new approval flow type.
5. In the **Name** box, enter a meaningful name.
6. Click **Save**.

Specify the approval flow type for a condition

Assign an approval flow type to a condition to be used in an approval flow. This task is performed in the Exception Monitor stage. Perform this procedure in the Exception Monitor settings in an Enterprise Designer dataflow.

1. In Enterprise Designer, double-click the Exception Monitor stage in a workflow to open its properties.
2. On the **Conditions** tab, click a condition and click the **Modify** button.
3. Under **Approval flow**, in the **Type** drop-down box, click an approval flow type.
4. Optional: To specify the data steward for the approval flow, click to choose a user in the **Assign to** box.

The default data steward for an approval flow, if you leave the **Assign to** setting blank, is the user who runs a job.

- Click **OK**, and click **OK** again.
- Click **File** > **Save** to save the changes to the workflow.

Enable approval flow types in security roles

Execute permission must be applied to an approval flow type within a security role before you can enable it in an approval flow. Perform this procedure in the Roles settings of the Management Console.

- Open the Spectrum Technology Platform **Management Console**.
On the Spectrum Technology Platform **Home** page, click **Platform Client tools** > **Web** > **Open Management Console**.
- On the **Management Console** page, click **System** > **Security**.
- On the **Security** page, click the **Roles** tab.
- Click the Add role button **+**.
- Click to expand the **Data Stewardship** section.
Approval flow types defined in the **Data Stewardship Settings** show in the table.
- For approval flow types that you want to include in approval flows, check the check box in the **Execute** column.
Execute permission is required to create an approval flow.

Assign approval flow role to new or existing users

To define reviewers in an approval flow, assign an approval flow role to new or existing users. Perform this procedure in the Users settings of the Management Console.

- Open the Spectrum Technology Platform **Management Console**.
On the Spectrum Technology Platform **Home** page, click **Platform Client tools** > **Web** > **Open Management Console**.
- On the **Management Console** page, click **System** > **Security**.
- On the **Security** page, click the **Users** tab.
- Open settings for a new or existing user.

To	Do this
Add a new user	Click the Add user + button to open the Add User page.
Edit an existing user	Click the check box next to a user and click the Edit user button  to open the Edit User page.

- For a new user, enter the user name and password.
- In the **Roles** box, check the check box next to the approval flow role that you want to use (see [Enable approval flow types in security roles](#) on page 219).
- Make sure that **Enabled** is set to **ON**.

8. Click the **Save** button.
9. Repeat steps **4** on page 219 through **8** on page 220 for each user that you want to include in the approval flow.

Define an approval flow

Define levels and specify the role and user for each level in the Management Console Approval Flow settings. Perform this procedure in the Approval Flows settings of the Management Console.

1. In the **Management Console**, click **Resources > Approval Flows**.
Approval flow types are listed in the **Type** column of the **Approval Flows** table.
2. In the **Type** column, click an approval flow type that you want to configure as an approval flow. This opens the Edit Approval flow page for the selected approval flow type.
3. Click the Add approval level button .
This adds a level to the approval flow. Each level will be associated with a reviewer defined in .
4. In the **Roles** column, select the role to be used for the level.
5. In the **User** column, select the user that records will be assigned to for this level.
If no user is selected, the record is assigned to the first user in alphabetical order.
6. In the **Days to Approve** column, select the maximum number of days in which the user must review and edit records.
7. To add additional levels to the approval flow, repeat steps **3** on page 220 through **7** on page 220.
8. Click to switch the **Enabled** box to **ON**.
9. Click the **Save** button.

Reviewing records in an approval flow

Perform steps in this task to edit, accept, or reject records in Tabular view.

1. On the **Data Stewardship Portal** Dashboard page, click the dataflow to view records that you want to review.
Alternatively, click **Editor**, then in the **Dataflow name** menu on the sidebar select the name of the dataflow.
2. Make sure that the records are displayed in Tabular view.
Click the Tabular view button  on the Editor toolbar.
3. In any record, click to edit a value.
Right-click the field to access **Cut**, **Copy**, and **Paste** commands. A field will be grayed out if it cannot be edited.
If you need to undo changes you made, select the record you want to undo and click the Undo changes  button.

After you have edited a field, you will notice a green triangle  appear in the upper-left corner of the field. This is a visual cue to remind you that the value of the field has been changed but is not yet saved.

4. Optional: You can add comments about your changes in the **Comment** box.
Comments are visible to other users and can be used to help keep track of the changes made to the record or to recommend additional changes. In an approval flow, you can describe additional changes necessary for approval.
5. Click the Save changes button .
This changes the record status to Edited.
6. After you have reviewed or edited values, you can click to Accept  or Reject  a record.
Clicking Accept  changes the record status to Accepted. Clicking Reject  changes the record status to Rejected.

Note: Anytime before you click the Save button , you can click Accept  to change the status from Accepted to Pending. This allows you to continue editing the record at a later time. Clicking Reject  changes the status to **Edited**

7. Click the Save changes button  button when you are finished editing records on a page.
Accepted records are assigned to the next approver in the approval flow. For the last approver in the approval flow, the record status is changed from Accepted to Resolved, and the records are made available for additional processing using the Read Exceptions stage.
Rejected records are assigned back to the data steward.

Resolving Duplicate Records

Duplicate resolution exceptions occur when the Spectrum Technology Platform cannot confidently determine whether a record is a duplicate of another. There are three ways to resolve duplicate records.

Note: Duplicate records can only be resolved with the Resolve Duplicates function on the Tabular view. However, you can still edit those records in the Form view.

Note: If you are working with Duplicate records in an Approval flow, see the section on [Approval Flows](#).

One approach is to group duplicate records together into collections. When you approve the records they can then be processed through a consolidation process to eliminate the duplicate records in each collection from your data.

Another approach is to edit the records so that they are more likely to be recognized as duplicates, for example correcting the spelling of a street name. When you approve the records, Spectrum Technology Platform reprocesses the records through a matching and consolidation process. If you corrected the records successfully, Spectrum Technology Platform will be able to identify the record as a duplicate.

Yet another approach to resolving duplicate records is to create a best of breed record. This combines the other two approaches by managing record collections and then editing one of the records in the collection to include fields from both the original and duplicate records. This record is then known as the best of breed record.

Making a Record a Duplicate of Another

Duplicate records are shown as groups of records in the Data Stewardship Portal. You can make a record a duplicate of another by moving it into the same group as the duplicate record.

To make a record a duplicate:

1. In the Exception Editor, select the record you want to work on then click the Resolve duplicates button .

The **Duplicate Resolution** view shows duplicate records. In this view, the initial record is highlighted in bold text. A column label is highlighted in blue text on any column of user data that contains values different from the value in the record. In these columns, values that do not exactly match the record are also highlighted in blue. Keep in mind that a mismatch may not be visible when a value is null or consists only of spaces.

For each record, the **MatchRecordType** column specifies one of three types of exceptions:

Suspect	A record that other records are compared to in order to determine if they are duplicates of each other. Each collection has one and only one suspect record.
Duplicate	A record that is a duplicate of the suspect record.
Unique	A record that has no duplicates.

2. If necessary, correct individual records as needed.
For more information, see [Editing and resolving records](#) on page 178.
3. In the **CollectionNumber** or **CandidateGroup** field, enter the number of the group that you want to move the record into. Alternatively, you can drag and drop records between groups.

The record is made a duplicate of the other records in the group.

In some cases you cannot move a record with a MatchRecordType value of "suspect" into another collection of duplicates.

Note: Records are grouped by either the CollectionNumber field or the CandidateGroup field depending the type of matching logic used in the dataflow that produced the

exceptions. Contact your Spectrum Technology Platform administrator if you would like additional information about matching.

- When you are done modifying the record, click Accept . This changes the record status to Accept. If you click Accept  again, the status is changed to Pending, so that the record can be reviewed later.
- To save your changes, click the Save changes  button. This changes the record status to Resolved and makes the record available for additional processing using the Read Exceptions stage.

Creating a New Group of Duplicate Records

In some situations you can create a new group of records that you want to make duplicates of each other. In other situations you cannot create new groups. Your ability to create new groups is determined by the type of Spectrum Technology Platform processing that generated the records.

- In the Editor, select the record you want to work on then click the Resolve duplicates button .

The **Duplicate Resolution** view shows duplicate records. In this view, the initial record is highlighted in bold text. A column label is highlighted in blue text on any column of user data that contains values different from the value in the record. In these columns, values that do not exactly match the record are also highlighted in blue. Keep in mind that a mismatch may not be visible when a value is null or consists only of spaces.

For each record, the **MatchRecordType** column specifies one of three types of exceptions:

Suspect	A record that other records are compared to in order to determine if they are duplicates of each other. Each collection has one and only one suspect record.
Duplicate	A record that is a duplicate of the suspect record.
Unique	A record that has no duplicates.

- If necessary, correct individual records as needed. For more information, see [Editing and resolving records](#) on page 178.
- Select a record that you want to put in the new collection then click the Create new collection button .

The new collection is automatically given a unique collection number, and the record you selected becomes a suspect record.

Note: If you do not see the **New Collection** button, you cannot create a new collection for the records you are working with. You can only create new collections if the dataflow that produced the exceptions contained an Interflow Match or an Intraflow Match stage, but not if it contained a Transactional Match stage. Contact your Spectrum Technology Platform administrator if you would like additional information about these matching stages.

4. Place additional records in the collection by entering the new collection's number in the record's `CollectionNumber` field.
5. When you are done modifying the record, click Accept .
This changes the record status to Accept. If you click Accept  again, the status is changed to Pending, so that the record can be reviewed later.
6. To save your changes, click the Save changes  button.
This changes the record status to Resolved and makes the record available for additional processing using the Read Exceptions stage.

Making a Record Unique

Follow steps in this procedure to change a record from a duplicate to a unique.

1. In the **MatchRecordType** field, click "Unique".
2. When you are done modifying the record, click Accept .
This changes the record status to Accept. If you click Accept  again, the status is changed to Pending, so that the record can be reviewed later.
3. To save your changes, click the Save changes  button.
This changes the record status to Resolved and makes the record available for additional processing using the Read Exceptions stage.

Fields Automatically Adjusted During Duplicate Resolution

When you modify records in the duplicate resolution view of the Data Stewardship Portal, some fields are automatically adjusted to reflect the record's new disposition.

Table 10: Records Processed by Interflow or Intraflow Match

Action	Values Automatically Applied to Fields
Moving a record from one collection to another	<p>If you move a record into a collection of duplicates:</p> <ul style="list-style-type: none"> • MatchRecordType: Duplicate • MatchScore: 100 • HasDuplicates: D (This field is only present if the dataflow contained an Interflow Match stage.) <p>If you move a duplicate record into the collection of unique records (collection 0):</p> <ul style="list-style-type: none"> • MatchRecordType: Unique • MatchScore: No change • HasDuplicates: U (This field is only present if the dataflow contained an Interflow Match stage.) <p>If you move a suspect record into the collection of unique records (collection 0):</p> <ul style="list-style-type: none"> • MatchRecordType: Unique • MatchScore: 0 • HasDuplicates: N (This field is only present if the dataflow contained an Interflow Match stage.)
Creating a new collection	<ul style="list-style-type: none"> • MatchRecordType: Suspect • MatchScore: No value • HasDuplicates: Y (This field is only present if the dataflow contained an Interflow Match stage.) <p>Note: If the record came from a dataflow that contained an Interflow Match stage only records with a value of "input_port_0" in the InterflowSourceType field can be a suspect record.</p>

Table 11: Records Processed by Transactional Match

Action	Values Automatically Applied to Fields
Change MatchRecordType to Duplicate	<ul style="list-style-type: none"> • HasDuplicates: D • MatchScore: 100
Change MatchRecordType to Unique	<ul style="list-style-type: none"> • HasDuplicates: U • MatchScore: unchanged

Action	Values Automatically Applied to Fields
Change HasDuplicates to D	<ul style="list-style-type: none"> • MatchRecordType: Duplicate • MatchScore: 100
Change HasDuplicates to U	<ul style="list-style-type: none"> • MatchRecordType: Unique • MatchScore: unchanged
Change HasDuplicates to Y	<ul style="list-style-type: none"> • MatchRecordType: Suspect • MatchScore: blank
Change HasDuplicates to N	<ul style="list-style-type: none"> • MatchRecordType: Suspect • MatchScore: blank

Creating a Best of Breed Record

Follow steps in this procedure to create a best of breed record in a collection.

Best of breed records can only be added to non-zero collections. Also, a best of breed record must be created for each collection within a group before you can accept all the records within that group. If a Best of Breed record does not exist for a collection and the user clicks save, the record status changes to Failed.

1. Select a record in the collection and click the Add Best of Breed Record button ).

This will make a copy of the record you selected and place it just after that record in the Data Stewardship Portal **Editor**. When you save the best of breed record, it will move to the bottom of the record collection. If you look at the best of breed record's history, you will see a link to the record from which it was copied.

Note: You can undo this step by selecting the best of breed record and then clicking the Remove Best of Breed Record button .

2. To combine fields into the best of breed record, drag and drop fields from the duplicate records into the best of breed record.

There are some limitations to drag-and-drop in this environment:

- The cell you are dragging from one record to another cannot be in an "edit" state (though the cell where you are dropping it can be).
- When you drop one cell onto another cell (whether in the same row or a different row) the second cell will take on the value of the first cell.
- You cannot drag and drop Boolean-type cells.

- You can drag non-string-type cells to string-type cells (such as date fields), but you cannot drop string-type cells to non-string-type cells.
- You cannot drop a cell onto another cell that uses lookup values (such as combination box editors), but you can drag a value from these cells to other string-type cells.
- You cannot drop a cell onto a read-only cell.

3. Click the Save changes  button.

Related reference

See [Write Exceptions options for more information on best of breed functionality](#). on page 324

Using Search Tools

The Data Stewardship Portal **Editor** provides search tools to assist you in looking up information that may help you edit and approve records and rerun them successfully. The tools include Maps, OpenNameParser, and other services you have licensed in Spectrum Technology Platform. Tools with icons that have a plus sign in the middle are premium tools and will incur charges when used.

Note: These search tools must first be configured on the **Data Stewardship Settings** page of Management Console. If none are selected on that page, none will appear on this page. Also, regardless of which tools are selected in Management Console, the services that appear to the user here are limited to those for which he has view permissions.

Using Maps

Follow steps in this procedure to search for map locations that correspond to records on the Data Stewardship Portal **Editor** page.

Before you can use this tool, the **Maps** service must be selected on the **Data Stewardship Settings** page of Management Console.

The Maps search tool displays the location of an address on a map. You can move the coordinates to different addresses, zoom in or out, and view adjacent regions.

Note: The Maps search tool is provided by Precisely LI GeoMap and GeoCode API services. You must be connected to the Internet to use this tool.

1. On the Data Stewardship Portal **Editor** page, click the record you want to research.
2. Below the records table, click **Search Tools** to expand the view.
3. In the **Service** field, select **Maps**.
4. To search for a location, do either of the following:
 - Enter an address in the **Input** box.
 - Select check boxes next to the **Field Name** column to search using field values.

For example, to search for an address on a map, you would select AddressLine1 and City. If you only want to view a city on a map, you would select City and StateProvince.

5. Click search button .

The result is displayed in the on the map in the form of a pin . Click the Zoom in  and Zoom out  buttons to adjust the detail. Drag the map to view adjacent regions. The address and the longitude and latitude are displayed in the address box above the map.



6. Drag the pin  to view the address and coordinates for another location on the map.

After completing an initial map search, you can click another record and the search button , and the map will update accordingly.

Using Spectrum Service Search Tools

To perform this procedure, you must first configure search tools on the **Data Stewardship Settings** page of Management Console. Otherwise **Search Tools** will not appear on the Data Stewardship Portal **Editor** page.

Precisely service search tools include all services for which you are licensed, such as ValidateAddress, GetPostalCodes, and so on. You can use these services within the Data Stewardship Portal **Editor** to look up and validate exception data that you are attempting to correct or augment.

Note that when using this feature you will only see services if you have view permissions for Services under the Platform group for role security. Likewise, in order to run services, you will need execute permissions for Services. However, these permissions can be modified by using Secured Entity Overrides. Using a combination of top-level permissions and overrides the administrator can manage the list of services that a particular user or role has access to in the Data Stewardship Portal **Service** box.

1. Select the record that contains the data you want to look up.
2. Below the records table, click to expand **Search Tools**.
3. In the **Service** box, choose the service you want to use, such as ValidateAddress or GetCandidateAddresses.
4. If the exception record contains fields used in that service but under different names, map the Service Input fields to the Exception Fields on the **Service Fields** tab.
For example, if you are using ValidateAddress and a record does not include an AddressLine1 field but instead includes an AddrLine1 field, select "AddrLine1" in the Exception Field column of the AddressLine1 row. You must have at least one input field mapped before running the service.

Note: The Data Stewardship Portal remembers the maps you create from service fields to exception fields as long as you are mapping records with the same field names, stage names, and dataflow names. For instance, if a record has a field named "AddrLine1" and you map it to "AddressLine1", it will remember this map as long as you are working with

records that contain "AddrLine1" and that were created in the same stage by the same dataflow.

5. Optional: Repeat step 4 on page 228 for Service Output fields.
This step is optional, but you must have at least one output field mapped before you can apply the service data. Note that you cannot map a read-only Exception Field to the Service Output fields.
6. Click the **Options** tab to select database resources and to view and change service options that were set in Management Console.
If you don't know the purpose of a particular option, hover over that option to see its description. Changes you make here will persist when used by the same user, dataflow, and stage of the record.
Note: If the service you are using requires a database, you must have first configured the database resource in Management Console. For example, if you are reviewing U.S. records using Validate Address, you must have configured a U.S. database in Management Console.
7. Click **Run Service**.
The updated record will appear on the **Result** tab, along with a status code indicating the success of the record. Fields that are mapped to records will be designated by an asterisk.
8. Select the result record and click **Apply Service Data** to transfer the data to the mapped fields of the record.
9. If you want to reprocess the updated record, click the **Approved** check box for that record and then click **Save**.

The Manage page

The **Manage** page on the Data Stewardship Portal enables a user with view and modify permissions to review and manage record activity for all assignees. It also provides the ability to reassign records from one user to another. If you have delete permissions, you can delete an entire group of records from the system based on dataflow name and job ID.

Reviewing Record Activity

In the Data Stewardship Portal, click **Manage** to display the **Manage** page.

Click the **Status** button to view the **Status** section of the **Manage** page. This section shows record activity by assignment or dataflow name. You can specify which one displays by clicking the

Assignments or **Dataflows** button near the top-right corner. **Assignments** provides the number of records assigned to each user as well as how many of those records have been approved and how many remain. It also shows the number of normally processed records that exist within groups that also contain records that have exceptions. Note that non-exception data will appear only if there are normally processed records within a group also containing records with exceptions and if there are no security overrides preventing this data from being shown.

Dataflows provides the number and percentage of records that have been approved for each dataflow. For dataflows containing more than one Exception Monitor stage, you can further break down this information by those stages. For example, if a dataflow contains two Exception Monitor stages, "NameExceptions" might be responsible for 37% of all exceptions in the dataflow and "DataExceptions" might be responsible for the other 63%. In this situation, the status of each **dataflow** is still calculated as a whole as it relates to all dataflows in the repository, while the status of each **stage** is calculated based on the total number of exceptions for the selected dataflow.

Additionally, you can filter the information that displays by entering search criteria in the Filter box. The list will dynamically auto-populate with dataflows or assignees that match the letters you type.

Click the **Progress** button to display the **Progress** section. This section shows the cumulative progress for all users or all dataflows whose individual statuses are shown on the Dashboard. Click **Assignments** at the top of the screen to view user progress. Click **Dataflows** to view dataflow progress. Use the scale to choose the metric for the scale (day, week, month) and the number of units to display.

Assigning Records

The **Maintenance > Assignment** section under on the **All Records** page enables you to reassign records from one user to another.

1. On the Data Stewardship Portal page, click **Manage**.
This displays the **Manage** page.
2. Click **Assignments** to expand the assignment options.
3. In the **User** field, select a user whose exceptions you want to assign to another user.
4. To reassign all records belonging to a user, skip to step 6 on page 233. To reassign a portion of a user's records, complete one or more of these fields:

Field	Description
Dataflow name	The name of the dataflow producing the records.
Stage label	The name of the stage producing the records.
Job ID	The ID assigned to the job containing the records.

Field	Description
Data domain	The kind of data assigned in the Exception Monitor.
Quality metrics	The kind of metric assigned in the Exception Monitor.
From date/time	The start date in a range of dates in which the records were created.
To date/time	The end date in a range of dates in which the records were created.
Record status	Select one or more status from the list. If none is selected, all records matching other filters are reassigned. For example, select status of Edited and Rejected records and only records with those status values are reassigned.

5. After completing your selections in the **User**, **Dataflow name**, and **Stage label** fields (at minimum), you can further refine the filter based on exception field values.
- Click the Add expression button **+**.
 - In the **Field Name** column, select the field you want to filter on.
 - In the **Operator** column, select one of the following:

Operator	Description
Is equal to	Looks for records that have exactly the value you specify. This can be a numeric value or a text value. For example, you can search for records with a MatchScore value of exactly 82, or records with a LastName value of "Smith."
Is not equal to	Looks for records that have any value other than the one you specify. This can be a numeric value or a text value. For example, you can search for records with any MatchScore value except 100, or records with any LastName except "Smith."
Is greater than	Looks for records that have a numeric value that is greater than the value you specify.
Is greater than or equal to	Looks for records that have a numeric value that is greater than or equal to the value you specify. For example, if you specify 50, you would see records with a value of 50 or greater in the selected field.

Operator	Description
Is less than	Looks for records that have a numeric value that is less than the value you specify.
Is less than or equal to	Looks for records that have a numeric value that is less than or equal to the value you specify. For example, if you specify 50, you would see records with a value of 50 or less in the selected field.
Contains	Looks for records that contain the value you specify in any position within the selected field. For example, if you filter for "South" in the AddressLine1 field, you would see records with "12 South Ave.", "9889 Southport St.", "600 South Shore Dr.", and "4089 5th St. South".
Does not contain	Looks for records that do not contain the value you specify in any position within the selected field. For example, if you filter for "South" in the AddressLine1 field but select "does not contain," you would not see records with "12 South Ave.", "9889 Southport St.", "600 South Shore Dr.", and "4089 5th St. South".
Starts with	Looks for records that start with a particular value in the selected field. For example, if you filter for "Van" in the LastName field you would see records with "Van Buren", "Vandenburg", or "Van Dyck".
Ends with	Looks for records that end with a particular value in the selected field. For example, if you filter for records that end with "burg" in the City field, you would see records with "Gettysburg", "Fredricksburg", and "Blacksburg".

- d) In the **Value** column, enter the value to use as the filtering criteria.

Note: The search value is case-sensitive. This means that searching for SMITH will return only records with "SMITH" in all upper case, but not "smith" or "Smith."

- e) To filter on more than one field, add multiple filters by clicking the add field filter icon again. For example, if you want all records with a LastName value of "SMITH" and a State value of "NY" you could use two filters, one for the LastName field and one for the State field.

This example would return all records with a value of "FL" in the State field:

The screenshot shows a filter configuration interface. At the top left, there are icons for adding (+), deleting (trash), and refreshing (refresh). Below these icons is a table with three columns: 'Field name', 'Operator', and 'Value'. The first row of the table has 'State' in the 'Field name' column, 'is equal to' in the 'Operator' column, and 'FL' in the 'Value' column. Each cell in the table has a small downward arrow icon, indicating a dropdown menu.

This example would return all records that do not have a PostalCode value of 60510:

Field name	Operator	Value
StateProvince	Is equal to	NY
PostalCode	Is not equal to	14226

This example would return all records with a StateProvince of "NY" with all postal codes except 14226.

Field name	Operator	Value
PostalCode	Is not equal to	60510

6. Click **Reassign**.
7. Optional: Select the number of exceptions you want to reassign.
You can assign all or some of the resulting exceptions to the new user. For example, if you enter "10" as the limit, only the first 10 records meeting the criteria will be reassigned and the remaining records meeting the criteria will not be reassigned.
8. Select another user in the **Reassign** box.
9. Click **Confirm**.

Deleting Records

Occasionally you may want to delete records from the exception repository. For instance, you could have residual records from testing the system or records that were mistakenly considered exceptions after processing, or you may want to process and delete approved records first and then re-run the same job again. The **Purge** section, under **Maintenance**, on the **All Records** page enables you to do this.

Note: You can also use the `bsm delete exceptions` command in the Command Line Interface (CLI) to delete records. For details about the CLI commands, see the Administration Utility section of the *Spectrum Technology Platform Administration Guide*.

1. In the Data Stewardship Portal, click **Manage**.
This displays the **Manage** page.
2. Under **Maintenance**, expand the **Purge** section.
3. In the **Dataflow name** box, select a dataflow.
4. In the **Job ID** box, either select a job ID or **All**
Select **All** to remove records from every job in the selected workflow.
5. Optional: To remove all performance data, you can check the **Remove data quality report data** check box.

If you do not check this check box, the records will be removed from the repository but performance data will still appear on the Data Quality page. This option is unavailable if you choose to add filter criteria in step 6 on page 234.

6. Optional: You can add filter criteria to specify records to be deleted.
 - a) In the **Stage label** box, select the stage that produced the records.
 - b) In the **Approval status** box, select one or more status from the list.
If none are selected, all records matching other filters are deleted.
 - c) For each filter you want to apply, click the Add expression button **+**.
In each case, specify a **Field name**, an **Operator**, and **Value** to filter based on a field values. You must have specified a **Job ID** and **Stage Label** in order to define field expressions in this step.
7. Click **Remove**.

The Data Quality Page

The Data Stewardship Portal **Data Quality** page provides information on trends among records.

Note: This page will be deactivated if you turn off data quality reporting on the Data Stewardship Settings page of Management Console.

Identifying Trends

The **Data Quality Trends** page depicts the following statistical information by dataflow and domain:

- Total number of conditions processed
- Total number of exception conditions
- Percentage of conditions that were processed successfully
- The trend of your data in 30-day intervals

This information can be further broken down by stage if you select a dataflow or by metric if you select a domain. The values that appear here are determined by dataflows with Exception Monitor stages that have been run over the last 30 days as well as their respective domains and metrics.

- If you are viewing information by domain, select a **Dataflow name** if you want to view information for a specific dataflow and a **Stage label** if you want to view information for a specific stage.

Note: You must select a single dataflow if you want to also filter the results based on a stage. Otherwise, you will see data for all dataflows.

- If you are viewing information by dataflow, select a **Domain** if you want to view information for a specific domain and a **Metric** if you want to view information for a specific metric. Otherwise, you will see data for all domains.
- Select a duration for the **Scale** to specify how far back you want the data to go. You can select any number of days, weeks, or months.
- Select a dataflow if you want to break out results by stage.
- Select a domain if you want to break out results by data quality metric.

As you make selections to view different parts of the data, you will notice that charts appear reflecting the status of that particular data. For instance, if you look at domain data, you could narrow that data down by a specific domain (such as "Address") and then by a particular metric (such as "Completeness"). The charts will update accordingly as you make different selections.

How to Video - Approval Flow in Data Stewardship

This video explains how you can build Approval Flow in Data Stewardship.

<https://www.youtube.com/embed/ymQTa35L1QQ>

Data Stewardship Settings

Introduction

Data Stewardship Settings provides the following tools for users with write permissions:

Tool	Description
Lookups	The Lookups tool provides a way for you to select from a list of values for a specific field when updating records in the Data Stewardship Portal Editor .
Domains	Domains specify the kind of data being evaluated.

Tool	Description
Metrics	Metrics specify the way in which data is measured.
Notifications	The Notifications feature enables you to have the system send a message to one or more email addresses when a designated number of exceptions are tied to a specific domain or metric.
Data Quality Reporting	The Data Quality Reporting settings configure preferences for tracking pass/fail conditions and KPIs.
Search Tools Services	Search Tools Services provides preferences for search tools in the Data Stewardship Portal Editor .
Options	These settings provide preferences for audit logs and progress tracking.
Approval Flow Types	The Approval Flow Types setting allows you to define types that associate records to an approval flow.

Accessing Data Stewardship Settings

Complete this procedure to access to access the **Data Stewardship Settings** page.

The **Data Stewardship Settings** page is located in the Management Console application.

1. In a web browser, navigate to this URL:

```
http://server:port/managementconsole
```

Where *server* is the server name or IP address of your Spectrum Technology Platform server and *port* is the HTTP port used by Spectrum Technology Platform. By default, the HTTP port is 8080.

2. Enter a valid user name and password.
3. Click the **Resources > Data Stewardship Settings**.

Lookups

The Lookups tool provides a way for you to select from a list of values for a specific field when updating records on the Data Stewardship Portal **Editor** page. This feature is particularly useful when you have several records with data in the same field that you want to change.

For example, you could have a set of exception records that contain banking data. One of the fields in that data could consist of codes that represent what kind of account is tied to the record (1=checking, 2=savings, 3=money market, and so on). whose addresses include ISO codes instead of names in the Country field, making those addresses unable to be validated. To correct this, you could create a lookup that provides ISO codes with their respective country names and make the corrections in the Data Stewardship Portal **Editor**, where you select the ISO code from a list that then populates the field with the country name tied to that ISO code.

Another benefit of using this tool is that it limits the options available for corrections, which reduces the possibility of further error. Using the same example of country names being incorrect or missing, by creating a lookup that provides a list of country names instead of requiring those names to be manually entered for each exception record, you ensure that those names are spelled correctly and are more likely to be validated when they are reprocessed.

What is the Lookup Process?

The lookup process involves three steps after you have reviewed exceptions and identified a recurring issue among those exceptions (such as invalid data in a country field):

- Create the lookup using values and/or labels of accurate data that will overwrite the bad data.
- Using the Write Exceptions stage in the dataflow that is producing the exception records, point the problematic field to the lookup you created and rerun the dataflow.
- Correct the exception records on the Data Stewardship Portal **Editor** page by overwriting bad data in the problematic field with good data from the lookup.

Creating Lookups

A lookup is made up of values or value/label pairs that contain data to replace existing, problematic data in a dataflow that is producing exception records. The value is what will replace the problematic data, and the label is what is displayed in a list you select from when using the lookup table to correct records on the Data Stewardship Portal **Editor** page.

Note: If you include only values in your lookup table, the values will also be used as labels.

You can populate a lookup by manually entering the information or by copying it from an external source and pasting it into the **Add many** dialog box. The external source can be a spreadsheet, a text file, or virtually any other file as long as the information is presented in one or two columns with either a comma, tab, or semicolon delimiter.

Note: When you use the **Add many** function and then click **Save**, any previously existing values or value/label pairs for that lookup will be deleted. However, after you have used the Add many function, you can manually add additional values or value/label pairs.

1. On the Management Console **Resources > Data Stewardship Settings** page, click **Lookups** in the sidebar.
2. Click the Add lookup button **+**.
3. Enter a name for the new lookup in the text box.

4. Add a value/label pair.

To manually add single value/label pairs:

- a) Click the Add lookup value button **+**.
- b) Enter **Value** and/or **Label** for the lookup pair.

To add lists of value/label pairs:

- a) Click the Add many button **+** to open the **Add many** dialog box.
- b) Configure **First column**, **Separator**, and **Second column** according to your list..
If you are pasting data from Microsoft Excel, use the Tab separator. If you select the wrong separator, the tool will import the entire line including the separator as the value or label (according to the **First column** selection).
- c) Type or paste in rows of values, separators, and labels.
If a row has no entry in the first column, the separator must still precede the second column entry.

After all value/label pairs have been added you can sort them in ascending or descending order on either the **Value** or the **Label** column.

Note: Once you have sorted the list, you can only change the order. You cannot return to the original in which values were initially added.

5. Repeat step **4** on page 238 to add additional value/label pairs.
6. Click **Save**.

Assigning Lookups

After creating a lookup, you need to assign that lookup to the field with problematic data in the Write Exceptions stage of the dataflow.

1. In Enterprise Designer, open the dataflow that is producing the exception records.
2. Open the Write Exceptions stage.
3. In the **Lookup name** column for the field with problematic data, select the lookup that contains the new, accurate data from the drop-down list and click **OK**.
4. Save and rerun the dataflow.

Correcting Records

After creating a lookup and assigning that lookup to a field in the dataflow, you need to correct the exception records in Data Stewardship.

1. On the Data Stewardship **Editor** page, select the dataflow that is producing the exception records.
2. For the first problematic record, click the field that you assigned the lookup to.
3. Click the drop-down button in that field and select the correct label for that record.

Remember: This label is not necessarily the same as the value. For example, if you want your field to have a value of "California" you might click a label that says "CA".

- Repeat step 3 on page 238 for each problematic record.
- Save the changed exceptions.

Modifying or Deleting Lookups

- On the Management Console **Resources > Data Stewardship Settings** page, click **Lookups** in the sidebar.
- Check the box next to the appropriate lookup.
- Modify or delete the lookup.

Option	Description
To modify the lookup	Click the Edit lookup button  , modify lookup pairs as necessary, and click Save .
To delete the lookup	Click the Delete lookup button  .

- Click the **Save** button.

Domains

Domains specify the kind of data being evaluated. This is used for reporting purposes to show which types of exceptions occur in your data. For example, if the condition evaluates the success or failure of address validation, the data domain could be "Address"; if the condition evaluates the success or failure of a geocoding operation, the data domain could be "Spatial", and so forth.

Note: The domains you establish here will serve as default options both for Data Stewardship Settings and the Exception Monitor stage.

You can select one of the predefined domains listed below or specify your own domain by clicking the **Add item** button and completing the fields as necessary. You can also edit domains by selecting a domain, clicking the **Edit item** button, and making any necessary changes. You can also filter the list of domains shown by entering search data in the **Filter** field. The results will update dynamically.

- Account—The condition checks a business or organization name associated with a sales account.
- Address—The condition checks address data, such as a complete mailing address or a postal code.
- Asset—The condition checks data about the property of a company, such as physical property, real estate, human resources, or other assets.

- Date—The condition checks date data.
- Email—The condition checks email data.
- Financial—The condition checks data related to currency, securities, and so forth.
- Name—The condition checks personal name data, such as a first name or last name.
- Phone—The condition checks phone number data.
- Product—The condition checks data about materials, parts, merchandise, and so forth.
- Spatial—The condition checks point, polygon, or line data which represents a defined geographic feature, such as flood plains, coastal lines, houses, sales territories, and so forth.
- SSN—The condition checks U.S. Social Security Number data.
- Uncategorized—Choose this option if you do not want to categorize this condition.

Metrics

Metrics specify the way in which data is measured. This is used for reporting purposes to show which types of exceptions occur in your data. For example, if the condition is designed to evaluate the record's completeness (meaning, for example, that all addresses contain postal codes) then you could specify "Completeness" as the data quality metric.

Note: The metrics you establish here will serve as default options both for Data Stewardship Settings and the Exception Monitor stage.

You can select one of the predefined metrics listed below or specify your own metric by clicking the **Add item** button and completing the fields as necessary. You can also edit metrics by selecting a metric, clicking the **Edit item** button, and making any necessary changes. You can also filter the list of metrics shown by entering search data in the **Filter** field. The results will update dynamically.

- Accuracy—The condition measures whether the data could be verified against a trusted source. For example, if an address could not be verified using data from the postal authority, it could be considered to be an exception because it is not accurate.
- Completeness—The condition measures whether data is missing essential attributes. For example, an address that is missing the postal code, or an account that is missing a contact name.
- Consistency—The condition measures whether the data is consistent between multiple systems. For example if your customer data system uses gender codes of M and F, but the data you are processing has gender codes of 0 and 1, the data could be considered to have consistency problems.
- Interpretability—The condition measures whether data is correctly parsed into a data structure that can be interpreted by another system. For example, social security numbers should contain only numeric data. If the data contains letters, such as xxx-xx-xxxx, the data could be considered to have interpretability problems.
- Recency—The condition measures whether the data is up to date. For example, if an individual moves but the address you have in your system contains the person's old address, the data could be considered to have a recency problem.
- Uncategorized—Choose this option if you do not want to categorize this condition.

- **Uniqueness**—The condition measures whether there is duplicate data. If the dataflow could not consolidate duplicate data, the records could be considered to be an exception.

Notifications

The Notifications feature enables you to have the system send a message to one or more email addresses when a designated number of exceptions are tied to a specific domain or metric.

A notification email includes a link to the failed records in the Exception Editor of the Data Stewardship Portal, where users can manually enter the correct data. To stop sending notifications to a particular email address, remove that address from the list of recipients in the Send notification to line of the Edit domain page.

Note: A mail server must be configured before you can successfully use a notification from within Data Stewardship Settings.

Related tasks

[Configure Notifications](#) on page 241

This procedure steps through configuring notifications for domains or metrics.

Configuring a Mail Server

Configure Notifications

This procedure steps through configuring notifications for domains or metrics.

A mail server must be configured before you can successfully use a notification from within Data Stewardship Settings.

1. In **Management Console**, click **Resources > Data Stewardship Settings**.
2. Click either **Domains** or **Metrics**.
3. Select the check box next a domain or metric, and click the Edit item button .
4. In the **Send notification to** box, select user names from the drop-down list or enter new email addresses to which notifications should be sent.
Users are configured in Management Console.
5. In the **Number of exceptions to trigger notification box**, select the number of exception records that should trigger a notification.
6. In the **Subject** box, enter the text that should be sent as the subject of the notification.
7. In the **Message** box, enter the message that should appear in the body of the notification.

You can use the following variables in the message to relay important information about the exceptions:

- `${jobID}`—The ID number of the job that produced the exception records.
- `${jobName}`—The name of the job that generated the exception records.
- `${userName}`—The name of the user whose job that generated the exception records.

- `${stageLabel}`—The name of the dataflow stage that produced the exception records.
 - `${link}`—A link to the Editor page in the Data Stewardship Portal, showing records for a particular dataflow.
8. Check the **Send reminder** check box if you want to send a reminder notification, and select the number of days that should pass before the reminder is sent.
 9. In the **Send reminder to** box, select user names or enter an email addresses.
 10. In **Reminder Subject**, enter the text that should be sent as the subject of the reminder notification.
 11. In **Reminder Message**, enter the message that should appear in the body of the reminder notification.
The reminder uses an additional variable:
 - `${Count}`—The number of exceptions for the specified dataflow or stage that have yet to be resolved.
 12. Check the **Remind daily** check box if you want a reminder notification to be sent every day until the exceptions are resolved.

Related concepts

[Notifications](#) on page 241

The Notifications feature enables you to have the system send a message to one or more email addresses when a designated number of exceptions are tied to a specific domain or metric.

Related tasks

[Configuring a Mail Server](#)

Data Quality Reporting

The Data Quality Reporting settings configure preferences for tracking pass/fail conditions and KPIs.

1. Click **Data quality reporting** to track pass or fail conditions in the Exception Monitor stage.
If you turn off this option, the **Data Quality** page Data Stewardship will contain no data. Likewise, the "Report only" field in the Exception Monitor stage for all dataflows will be disabled.
2. In the **Retention** drop-down, select how long, in months, that data should be retained.

Configuring Key Performance Indicators

The **KPI Configuration** section of the Data Quality Reporting tab enables you to designate key performance indicators (KPIs) for your data and assign notifications for when those KPIs meet certain conditions.

1. Click **Add a KPI**.
2. Enter a **Name** for the key performance indicator. This name must be unique on your Spectrum Technology Platform server.

3. Select one of the data quality **Metrics** for the key performance indicator; if you do not make a selection, this key performance indicator will be tied to all **metrics**.
4. Select a **Dataflow name** for the key performance indicator; if you do not make a selection, this key performance indicator will be tied to all Data Stewardship dataflows.
5. Select a **Stage label** for the key performance indicator; if you do not make a selection, this key performance indicator will be tied to all Data Stewardship stages in your dataflows.
6. Select a data **Domain** for the key performance indicator; if you do not make a selection, this key performance indicator will be tied to all **domains**. Note that selecting a Domain here will cause the Condition field to be disabled.
7. Select a **Condition** for the key performance indicator. If you do not make a selection, this key performance indicator will default to "All". Note that to select a condition, you must first have selected "All" in the Domain field. Once a Condition has been selected, the Domain field will become disabled.
8. Select a **KPI period** to designate the intervals for which you want Data Stewardship to monitor your data and send notifications. For example, if you select "1" and "Monthly", a KPI notification will be sent when the percentage of exceptions has increased per the threshold or variance over a month-to-month period of time.
9. Provide a percentage for either a **Variance** or a **Threshold**. Variance values represent the increased percentage of failures in exception records since the last time period. Threshold values represent the percentage of failures at which you want the notifications to be sent. Its value must be 1 or greater.
10. Select email addresses from the list or enter email addresses for the **Recipients** who should be notified when these conditions are met. When possible, this field will auto-complete as you enter email addresses. You do not need to separate addresses with commas, semicolons, or any other punctuation.
11. Enter the **Subject** you want the notification email to use.
12. Enter the **Message** you want the notification to relay when these conditions are met.
13. Click **OK**. The new KPI will appear among any other existing KPIs. You can sort KPIs on any of the columns containing data.
14. Click **Save**.

You can modify and remove KPIs by selecting a KPI and clicking either **Edit selected KPI** or **Delete selected KPI**.

Search Tools Services

Search Tools Services provides preferences for search tools in the Data Stewardship Portal **Editor**.

1. Select which search tool services you want to be available in the Data Stewardship Portal **Editor**. The list of available services is based on user permissions and is populated by your licensed modules and services within Spectrum Technology Platform. Use the **Filter** to narrow the list of services based on filter criteria.

2. Click **Premium** to indicate to users that they will accrue additional fees when they use these services (such as Dun & Bradstreet services).

Options

These settings provide preferences for audit logs and progress tracking.

1. Click **Audit exception events** to have Data Stewardship maintain a log of when exception records are created, read, updated, or deleted.
2. Click **Track progress** to track when exception records are approved in Data Stewardship. If you turn this option off, progress charts will not appear on the Data Stewardship.

Approval Flow Types

The **Approval Flow Types** setting allows you to define types that associate records to an approval flow.

Approval flow types are used when defining conditions in the Exception Monitor stage. If a record meets the criteria of a condition, the record will then be associated to that approval flow. For more information, see [Approval Flows](#).

- To create a new approval flow type, click the Add type button .
- To delete existing approval flow types, check the check box next to type names that you want to delete, and click the Delete type button .

9 - Accessing Stages in Enterprise Designer

In this section

Advanced Matching.....	246
Data Stewardardship.....	310
Data Normalization.....	325
Universal Name.....	344



Advanced Matching

Spectrum Advanced Matching

Spectrum Advanced Matching matches records between and within any number of input files. You can also use Spectrum Advanced Matching to match on a variety of fields including name, address, name and address, or non-name and address fields, such as social security number or date of birth.

Components

Spectrum Advanced Matching consists of:

- **Best-of-Breed**—This stage selects a best-of-breed record from the duplicates cluster by selecting a template record, then using that record to build a composite record that becomes the survivor record.
- **Candidate Finder**—This stage finds candidate records that will form the set of potential matches that the Transactional Match stage will evaluate. Candidate Finder is used with Transactional Match.
- **Duplicate Synchronization**—This stage allows you to specify fields from a collection of records to copy (post) to the corresponding fields of all records in the collection.
- **Filter**—This stage allows you to specify the criteria that records must satisfy to be retained or removed from a collection of duplicate records, either for further processing downstream or for your output file.
- **Interflow Match**—This stage identifies matches between similar data records across multiple input streams.
- **Intraflow Match**—This stage identifies matches between similar data records within a single input stream.
- **Match Analysis**—Use Match Analysis to analyze and compare results from matching stages to understand the outcome of the matching process and improve overall matching results.
- **Match Key Generator**—This stage creates a non-unique key shared by all like records. You have the option of grouping records that share the same key for comparison.
- **Private Match**—This stage enables two entities to compare datasets and identify common records, or customers, without compromising sensitive information. Special encryption helps to retain security during matching, and a decryption function shows the output of the matched data.
- **Transactional Match**—This stage matches suspect transactions against a database using the Candidate Finder Stage to query and return potential candidate records.

- **Write to Search Index**—This stage enables you to create a full-text index based on the data that the stage is accepting for processing. You can retain data in a dedicated search index to improve response times when searching the Candidate Finder.

Best of Breed

Best of Breed consolidates duplicate records by selecting the best data in a duplicate record collection and creating a new consolidated record using the best data. This "super" record is known as the best of breed record. You define the rules to use in selecting records to process. When processing completes, the best of breed record is retained by the system.

Options

The following table lists the options for Best of Breed.

Option Name	Description / Valid Values
Group by	Specifies the field to use to create groups of records to merge into a single best of breed record, creating one best of breed record from each group. In cases where you have used a matching stage earlier in the dataflow, you should select the CollectionNumber field to use the collections created by the matching stage as the groups. However, if you want to group records by some other field, choose the field here. For example, if you want to merge all records that have the same value in the AccountNumber field into one best of breed record, you would select AccountNumber.
Sort	If you specify a field in the Group by field, check this box to sort the records by the value in the field you chose. This option is enabled by default.

Option Name	Description / Valid Values
Advanced	<p>Click this button to specify sort performance options. By default, the sort performance options specified in Management Console, which are the default performance options for your system, are in effect. If you want to override your system's default performance options, check the Override sort performance options box then specify the values you want in these fields:</p> <p>In memory record limit Specifies the maximum number of data rows a sorter will hold in memory before it starts paging to disk. By default, a sort of 10,000 records or less will be done in memory and a sort of more than 10,000 records will be performed as a disk sort. The maximum limit is 100,000 records. Typically an in-memory sort is much faster than a disk sort, so this value should be set high enough so that most of the sorts will be in-memory sorts and only large sets will be written to disk.</p> <p style="padding-left: 40px;">Note: Be careful in environments where there are jobs running concurrently because increasing the In memory record limit setting increases the likelihood of running out of memory.</p> <p>Maximum number of temporary files Specifies the maximum number of temporary files that may be used by a sort process. Using a larger number of temporary files can result in better performance. However, the optimal number is highly dependent on the configuration of the server running Spectrum Technology Platform. You should experiment with different settings, observing the effect on performance of using more or fewer temporary files. To calculate the approximate number of temporary files that may be needed, use this equation:</p> $\frac{(NumberOfRecords \times 2)}{InMemoryRecordLimit} = NumberOfTempFilesN$ <p style="padding-left: 40px;">Note: The maximum number of temporary files cannot be more than 1,000.</p> <p>Enable compression Specifies that temporary files are compressed when they are written to disk.</p> <p style="padding-left: 40px;">Note: The optimal sort performance settings depends on your server's hardware configuration. You can use this equation as a general guideline to produce good sort performance: $(InMemoryRecordLimit \times MaxNumberOfTempFiles \div 2) \geq TotalNumberOfRecords$</p>
Keep original records	<p>Select this option to retain all records in the collection along with the best of breed record. Clear the option if you want only the best of breed record.</p>
Use first record	<p>Select this option if you want Best of Breed to automatically select the first record in the collection as the template record. The template record is the record upon which the best of breed record is based.</p>

Option Name	Description / Valid Values
Define template record	Select this option to define rules for selecting the template record. For more information, see Defining Template Record Rules on page 249.

Defining Template Record Rules

In Best of Breed processing, the template record is the record in a collection that is used to create the best of breed record. The template record is used as the starting point for constructing the best of breed record and is modified based on the best of breed settings you define. The Best of Breed stage can select the template record automatically, or you can define rules for selecting the template record. This topic describes how to define rules for selecting the template record.

Template rules are written by specifying the field name, an operator, a value type, and a value. Here is an example of template record options:

Field Name: MatchScore
 Field Type: Numeric
 Operator: Equal
 Value Type: String
 Value: 100

This template rule selects the record in the collection where the Match Score is equal to the value of 100.

The following procedure describes how to define a template record rule in the Best of Breed stage.

1. In the Best of Breed stage, under **Template Record Settings**, select the option **Define template record**.
2. In the tree, click **Rules**.
3. Click **Add Rule**.
4. Complete the following fields.

Option	Description
Field name	Specifies the name of the dataflow field whose value you want to evaluate to determine if the record should be the template record.

Option	Description		
Field Type	Specifies the type of data in the field. One of the following:		
	<table border="0"> <tr> <td data-bbox="526 384 683 415">Non-Numeric</td> <td data-bbox="773 384 1414 443">Choose this option if the field contains non-numeric data (for example, string data).</td> </tr> </table>	Non-Numeric	Choose this option if the field contains non-numeric data (for example, string data).
	Non-Numeric	Choose this option if the field contains non-numeric data (for example, string data).	
<table border="0"> <tr> <td data-bbox="526 457 634 489">Numeric</td> <td data-bbox="773 457 1365 516">Choose this option if the field contains numeric data (for example, double, float, and so on).</td> </tr> </table>	Numeric	Choose this option if the field contains numeric data (for example, double, float, and so on).	
Numeric	Choose this option if the field contains numeric data (for example, double, float, and so on).		
Operator	Specifies the type of comparison you want to use to evaluate the field. One of the following:		
	<table border="0"> <tr> <td data-bbox="526 678 634 709">Contains</td> <td data-bbox="735 678 1414 737">Determines if the field contains the value specified. For example, "sailboat" contains the value "boat".</td> </tr> </table>	Contains	Determines if the field contains the value specified. For example, "sailboat" contains the value "boat".
	Contains	Determines if the field contains the value specified. For example, "sailboat" contains the value "boat".	
	<table border="0"> <tr> <td data-bbox="526 751 602 783">Equal</td> <td data-bbox="735 751 1338 783">Determines if the field contains the exact value specified.</td> </tr> </table>	Equal	Determines if the field contains the exact value specified.
	Equal	Determines if the field contains the exact value specified.	
	<table border="0"> <tr> <td data-bbox="526 798 683 829">Greater Than</td> <td data-bbox="735 798 1398 856">Determines if the field value is greater than the value specified. This operation only works on numeric fields.</td> </tr> </table>	Greater Than	Determines if the field value is greater than the value specified. This operation only works on numeric fields.
	Greater Than	Determines if the field value is greater than the value specified. This operation only works on numeric fields.	
	<table border="0"> <tr> <td data-bbox="526 871 716 930">Greater Than Or Equal To</td> <td data-bbox="735 871 1414 930">Determines if the field value is greater than or equal to the value specified. This operation only works on numeric fields.</td> </tr> </table>	Greater Than Or Equal To	Determines if the field value is greater than or equal to the value specified. This operation only works on numeric fields.
	Greater Than Or Equal To	Determines if the field value is greater than or equal to the value specified. This operation only works on numeric fields.	
	<table border="0"> <tr> <td data-bbox="526 945 618 976">Highest</td> <td data-bbox="735 945 1414 1129">Compares the field's value for all the records group and determines which record has the highest value in the field. For example, if the fields in the group contain values of 10, 20, 30, and 100, the record with the field value 100 would be selected. This operation only works on numeric fields. If multiple records are tied for the longest value, one record is selected.</td> </tr> </table>	Highest	Compares the field's value for all the records group and determines which record has the highest value in the field. For example, if the fields in the group contain values of 10, 20, 30, and 100, the record with the field value 100 would be selected. This operation only works on numeric fields. If multiple records are tied for the longest value, one record is selected.
	Highest	Compares the field's value for all the records group and determines which record has the highest value in the field. For example, if the fields in the group contain values of 10, 20, 30, and 100, the record with the field value 100 would be selected. This operation only works on numeric fields. If multiple records are tied for the longest value, one record is selected.	
	<table border="0"> <tr> <td data-bbox="526 1144 634 1176">Is Empty</td> <td data-bbox="735 1144 1166 1176">Determines if the field contains no value.</td> </tr> </table>	Is Empty	Determines if the field contains no value.
	Is Empty	Determines if the field contains no value.	
	<table border="0"> <tr> <td data-bbox="526 1190 683 1222">Is Not Empty</td> <td data-bbox="735 1190 1179 1222">Determines if the field contains any value.</td> </tr> </table>	Is Not Empty	Determines if the field contains any value.
Is Not Empty	Determines if the field contains any value.		
<table border="0"> <tr> <td data-bbox="526 1236 651 1268">Less Than</td> <td data-bbox="735 1236 1414 1295">Determines if the field value is less than the value specified. This operation only works on numeric fields.</td> </tr> </table>	Less Than	Determines if the field value is less than the value specified. This operation only works on numeric fields.	
Less Than	Determines if the field value is less than the value specified. This operation only works on numeric fields.		
<table border="0"> <tr> <td data-bbox="526 1310 683 1369">Less Than Or Equal To</td> <td data-bbox="735 1310 1382 1369">Determines if the field value is less than or equal to the value specified. This operation only works on numeric fields.</td> </tr> </table>	Less Than Or Equal To	Determines if the field value is less than or equal to the value specified. This operation only works on numeric fields.	
Less Than Or Equal To	Determines if the field value is less than or equal to the value specified. This operation only works on numeric fields.		
<table border="0"> <tr> <td data-bbox="526 1383 630 1415">Longest</td> <td data-bbox="735 1383 1414 1535">Compares the field's value for all the records group and determines which record has the longest (in bytes) value in the field. For example, if the group contains the values "Mike" and "Michael", the record with the value "Michael" would be selected. If multiple records are tied for the longest value, one record is selected.</td> </tr> </table>	Longest	Compares the field's value for all the records group and determines which record has the longest (in bytes) value in the field. For example, if the group contains the values "Mike" and "Michael", the record with the value "Michael" would be selected. If multiple records are tied for the longest value, one record is selected.	
Longest	Compares the field's value for all the records group and determines which record has the longest (in bytes) value in the field. For example, if the group contains the values "Mike" and "Michael", the record with the value "Michael" would be selected. If multiple records are tied for the longest value, one record is selected.		
<table border="0"> <tr> <td data-bbox="526 1549 618 1581">Lowest</td> <td data-bbox="735 1549 1414 1734">Compares the field's value for all the records group and determines which record has the lowest value in the field. For example, if the fields in the group contain values of 10, 20, 30, and 100, the record with the field value 10 would be selected. This operation only works on numeric fields. If multiple records are tied for the longest value, one record is selected.</td> </tr> </table>	Lowest	Compares the field's value for all the records group and determines which record has the lowest value in the field. For example, if the fields in the group contain values of 10, 20, 30, and 100, the record with the field value 10 would be selected. This operation only works on numeric fields. If multiple records are tied for the longest value, one record is selected.	
Lowest	Compares the field's value for all the records group and determines which record has the lowest value in the field. For example, if the fields in the group contain values of 10, 20, 30, and 100, the record with the field value 10 would be selected. This operation only works on numeric fields. If multiple records are tied for the longest value, one record is selected.		
<table border="0"> <tr> <td data-bbox="526 1749 699 1780">Most Common</td> <td data-bbox="735 1749 1414 1835">Determines if the field value contains the value that occurs most frequently in this field among the records in the group. If two or more values are most common, no action is taken.</td> </tr> </table>	Most Common	Determines if the field value contains the value that occurs most frequently in this field among the records in the group. If two or more values are most common, no action is taken.	
Most Common	Determines if the field value contains the value that occurs most frequently in this field among the records in the group. If two or more values are most common, no action is taken.		
<table border="0"> <tr> <td data-bbox="526 1850 646 1881">Not Equal</td> <td data-bbox="735 1850 1414 1881">Determines if the field value is not the same as the value specified.</td> </tr> </table>	Not Equal	Determines if the field value is not the same as the value specified.	
Not Equal	Determines if the field value is not the same as the value specified.		

Option	Description
Value type	<p>Specifies the type of value you want to compare to the field's value. One of the following:</p> <p>Note: This option is not available if you select the operator Highest, Lowest, or Longest.</p> <p>Field Choose this option if you want to compare another dataflow field's value to the field.</p> <p>String Choose this option if you want to compare the field to a specific value.</p>
Value	<p>Specifies the value to compare to the field's value. If you selected Field in the Field type field, select a dataflow field. If you selected String in the Value type field, type the value you want to use in the comparison.</p> <p>Note: This option is not available if you select the operator Highest, Lowest, or Longest.</p>

- Click **OK**.
- If you want to specify additional rules, click **Add Rule**.

If you add additional rules, you will have to select a logical operator to use between each rule. Choose **And** if you want the new rule and the previous rule to both pass in order for it to be selected as the template record. Select **Or** if you want either the previous rule or the new rule to pass in order for the record to be selected as the template record.

You have now configured rules to use to select the template record. Configure the best of breed settings to complete the configuration of the Best of Breed stage.

Defining Best of Breed Rules and Actions

Best of Breed rules and actions work together to determine which fields from duplicate records in a collection to copy to the Best of Breed record. Rules test values in a record and if the record passes the rules, the data is copied from the record to the template record. Actions define which data to copy, and which field in the template record should receive the data. After all the rules and actions are executed, the template record will be the best of breed record.

Rules and actions can be grouped together into conditions, and you can have multiple conditions. This allows you

- In the Best of Breed stage, under **Best of Breed Settings**, click the **Rules** node in the tree.
- Click **Add Rule**.
- Complete the following fields:

Option	Description
Field name	Specifies the name of the dataflow field whose value you want to evaluate to determine if the condition is met and the associated actions should be taken.
Field Type	Specifies the type of data in the field. One of the following: Non-Numeric Choose this option if the field contains non-numeric data (for example, string data). Numeric Choose this option if the field contains numeric data (for example, double, float, and so on).

Option	Description
Operator	Specifies the type of comparison you want to use to evaluate the field. One of the following:
Contains	Determines if the field contains the value specified. For example, "sailboat" contains the value "boat".
Equal	Determines if the field contains the exact value specified.
Greater Than	Determines if the field value is greater than the value specified. This operation only works on numeric fields.
Greater Than Or Equal To	Determines if the field value is greater than or equal to the value specified. This operation only works on numeric fields.
Highest	Compares the field's value for all the records group and determines which record has the highest value in the field. For example, if the fields in the group contain values of 10, 20, 30, and 100, the record with the field value 100 would be selected. This operation only works on numeric fields. If multiple records are tied for the longest value, one record is selected.
Is Empty	Determines if the field contains no value.
Is Not Empty	Determines if the field contains any value.
Less Than	Determines if the field value is less than the value specified. This operation only works on numeric fields.
Less Than Or Equal To	Determines if the field value is less than or equal to the value specified. This operation only works on numeric fields.
Longest	Compares the field's value for all the records group and determines which record has the longest (in bytes) value in the field. For example, if the group contains the values "Mike" and "Michael", the record with the value "Michael" would be selected. If multiple records are tied for the longest value, one record is selected.
Lowest	Compares the field's value for all the records group and determines which record has the lowest value in the field. For example, if the fields in the group contain values of 10, 20, 30, and 100, the record with the field value 10 would be selected. This operation only works on numeric fields. If multiple records are tied for the longest value, one record is selected.
Most Common	Determines if the field value contains the value that occurs most frequently in this field among the records in the group. If two or more values are most common, no action is taken.
Not Equal	Determines if the field value is not the same as the value specified.

Option	Description
Value type	<p>Specifies the type of value you want to compare to the field's value. One of the following:</p> <p>Note: This option is not available if you select the operator Highest, Lowest, or Longest.</p> <p>Field Choose this option if you want to compare another dataflow field's value to the field.</p> <p>String Choose this option if you want to compare the field to a specific value.</p>
Value	<p>Specifies the value to compare to the field's value. If you selected Field in the Field type field, select a dataflow field. If you selected String in the Value type field, type the value you want to use in the comparison.</p> <p>Note: This option is not available if you select the operator Highest, Lowest, or Longest.</p>

- Click **OK**.
- If you want to specify additional rules for this condition, click **Add Rule**.

If you add additional rules, you will have to select a logical operator to use between each rule. Choose **And** if you want the new rule and the previous rule to both pass in order for the condition to be met and the associated actions taken. Select **Or** if you want either the previous rule or the new rule to pass in order for the condition to be met.

- Click the **Actions** node in the tree.
- Click **Add Action**.
- Complete the following fields.

Option	Description
Source type	<p>Specifies the type of data to copy to the best of breed record. One of the following.</p> <p>Field Choose this option if you want to copy a value from a field to the best of breed record.</p> <p>String Choose this option if you want to copy a constant value to the best of breed record.</p>
Source data	<p>Specifies the data to copy to the best of breed record. If the source type is Field, select the field whose value you want to copy to the destination field. If the source type is String, specify a constant value to copy to the destination field.</p>

Option	Description
Destination	Specifies the field in the best of breed record to which you want to copy the data specified in the Source data field.
Accumulate source data	<p>If the data in the Source data field is numeric data, you can enable this option to combine the source data for all duplicate records and put the total value in the best of breed record.</p> <p>For example, if there were three duplicate records in the group and they contained these values in the Deposits field:</p> <p>100.00 20.00 5.00</p> <p>Then all three values would be combined and the total value, 125.00, would be put in the best of breed record's Deposits field.</p>

9. Click **OK**.
10. If you want to specify additional actions to take for this condition, click **Add Action** and repeat the above steps.
11. To add another condition, click the root condition in the tree then click **Add Condition**.

Example Best of Breed Rule and Action

This Best of Breed rule selects the record where the Match Score is equal to the value of 100. The Account Number data that corresponds to the selected field is then copied to the AccountNumber field on the Best of Breed record.

Rule
 Field Name: MatchScore
 Field Type: Numeric
 Operator: Equal
 Value Type: String
 Value: 100

Action
 Source Type: Field
 Source Data: AccountNumber
 Destination: AccountNumber

Output

Table 12: Best of Breed Output

Field Name	Format	Description / Valid Values
CollectionRecordType	String	<p>Identifies the template and Best of Breed records in a collection of duplicate records. The possible values are:</p> <p>Primary The record is the selected template record in a collection.</p> <p>Secondary The record is not the selected template record in a collection.</p> <p>BestOfBreed The record is the newly created best of breed record in the collection.</p> <p>Note: The Primary and Secondary values are generated only when you select the Define template record option in the Best of Breed Options window.</p>

Candidate Finder

Candidate Finder obtains the candidate records that will form the set of potential matches. Database searches work in conjunction with Transactional Match, and Search Index searches work independently from Transactional Match. Depending on the format of your data, Candidate Finder may also need to parse the name or address of the suspect record, the candidate records, or both.

Candidate Finder also enables full-text index searches and helps in defining both simple and complex search criteria against characters and text using various search types (Any Word Starts With, Contains, Contains All, Contains Any, Contains None, Fuzzy, Pattern, Proximity, Range, Wildcard) and conditions (All True, Any True, None True).

Database Options

The Candidate Finder dialog enables you to define SQL statements that retrieve potential match candidates from a database, as well as map the columns that you select from the database to the field names that are defined in your dataflow.

Table 13: Candidate Finder Database Options

Option Name	Description / Valid Values
Finder type	Select Database.
Connection	Select the database that contains the candidate records. You can select any connection configured in Management Console. To connect to a database not listed, configure a connection to that database in Management Console, then close and reopen Candidate Finder to refresh the connection list. Note: The Dataflow Options feature in Enterprise Designer enables the connection name to be exposed for configuration at runtime.
SQL statement	Type a SQL statement in the text box as described in Defining the SQL Query on page 257
Field Map tab	Choose field mapping settings as described in Mapping Database Columns to Stage Fields on page 258.
Preview tab	Click this tab to enter a sample match key to test your SQL <code>SELECT</code> statement or your index query.

Defining the SQL Query

You can type any valid SQL select statement into the text box on the **Candidate Finder Options** dialog.

Note: `Select *` is not valid.

For example, assume you have a table in your database called `Customer_Table` that has the following columns:

- `Customer_Table`
- `Cust_Name`
- `Cust_Address`
- `Cust_City`
- `Cust_State`
- `Cust_Zip`

To retrieve all the rows from the database, you might construct a query similar to the following:

```
SELECT Cust_Name, Cust_Address, Cust_City, Cust_State, Cust_Zip from
Customer_Table;
```

You will rarely want to match your transaction against all the rows in the database. To return only relevant candidate records, add a `WHERE` clause using variable substitution. Variable substitution refers to a special notation that you will use to cause the Candidate Selection engine to replace the variable with the actual data from your suspect record.

To use variable substitution, enclose the field name in braces preceded by a dollar sign using the form `${FieldName}`. For example, the following query will return only those records that have a value in `Cust_Zip` that matches the value in `PostalCode` on the suspect record.

```
SELECT Cust_Name, Cust_Address, Cust_City, Cust_State, Cust_Zip
FROM Customer_Table
WHERE Cust_Zip = ${PostalCode};
```

For SQL 2000, the data type needs to be identical to the data type for Candidate Finder. The JDBC driver sets the Candidate Finder input variable (Ex: `${MatchKey}`) that is used in the `WHERE` clause to a data type of `nVarChar(4000)`. If the data in the database is set to a data type of `VarChar`, SQL Server will ignore the index on the database. If the index is ignored, then performance will be degraded. Therefore, use the following query for SQL 2000:

```
SELECT Cust_Name, Cust_Address, Cust_City, Cust_State, Cust_Zip
FROM Customer_Table
WHERE Cust_Zip = CAST(${PostalCode} AS VARCHAR(255));
```

Mapping Database Columns to Stage Fields

If the column names in your database match the Component Field names exactly, they are automatically mapped to the corresponding Stage Fields. If they are not named exactly the same, you will need to use the Selected Fields (columns from the database) to map to the Stage Fields (field names defined in the dataflow).

For example, consider a table named `Customer_Table` with the following columns:

- `Cust_Name`
- `Cust_Address`
- `Cust_City`
- `Cust_State`
- `Cust_Zip`

When you retrieve these records from the database, you need to map the column names to the field names that are used by Transactional Match and other components in your dataflow. For example, `Cust_Address` might be mapped to `AddressLine1`, and `Cust_Zip` would be mapped to `PostalCode`.

1. Select the drop-down list under **Selected Fields** in the **Candidate Finder Options** dialog. Then, select the database column `Cust_Zip`.

2. Select the drop-down list under **Stage Fields**. Then, select the field to which you want to map.

For example, if you want to map Cust_Zip to Postal Code, first select Cust_Zip under Selected fields and then select PostalCode on the corresponding Stage Field row.

Alternate Method for Mapping Fields

You can use special notation in your SQL query to perform the mapping. To do this, enclose the field name you want to map to in braces after the column name in your query. When you do this, the selected fields are automatically mapped to the corresponding stage fields.

For example,

```
select Cust_Name {Name}, Cust_Address {AddressLine1},
       Cust_City {City}, Cust_State {StateProvince},
       Cust_Zip {PostalCode}
from Customer
where Cust_Zip = ${PostalCode};
```

Configuring the Connection Name at Runtime

The Connection name can be configured and passed at runtime if it is exposed as a dataflow option. This enables you to run your dataflow while using a different connection name.

1. In Enterprise Designer, open a dataflow that uses the Candidate Finder stage.
2. Save and expose that dataflow.
3. Go to **Edit > Dataflow Options**.
4. In the **Map dataflow options to stages** table, expand Candidate Finder and edit options as necessary. Check the box for the option you want to edit, then change the value in the **Default value** drop-down.
5. Optional: Change the name of the options in the **Option label** field.
6. Click **OK** twice.

Search Index Options

The Candidate Finder dialog enables you to conduct a simple search that matches input field values within search indexes or an advanced search to build matching rules that retrieve potential match candidates from search indexes.

Simple Search Index Options

Table 14: Candidate Finder Options

Option Name	Description / Valid Values
Finder type	Select Search Index.
Name	Select the appropriate index that was created using the Write to Search Index on page 305 stage under the Advanced Matching deployed stages in Enterprise Designer.
Starting record	Enter the record number on which search results should begin. The default is 1.
Maximum results	<p>Enter the maximum number of results you want the search index to return. Default is 10.</p> <p>Note: If the maximum results is arbitrarily large, process those in batches, using the Fetch Batch Size field.</p>
Fetch Batch Size	<p>If the Maximum results is arbitrarily large, enter the size of batches in which you want the results to be processed. This optimizes processing of large number of records. Default is 10000.</p> <p>The recommended Fetch Batch Size is a value lesser than Maximum results and if the Fetch Batch Size is greater than Maximum results, the records are processed in a single batch.</p> <p>Note: This field is applicable only to cluster supported search engine and not to the legacy search engine.</p>
Return match count	<p>Returns the total number of matches that were made. For example, if you use the default of "10" for the Maximum results field above, only 10 results will be returned. However, if you check this box, the TotalMatchCount output field will tell you how many matches were made during processing.</p>
Index search type	Determines the type of index search you want to conduct. Select Simple search .
Index Fields	Select the index field(s) you want to use for comparison in the simple search.

Option Name	Description / Valid Values
Input field	Select the input field you want to use for comparison in the simple search.

Option Name	Description / Valid Values
-------------	----------------------------

Input analyzer	
----------------	--

Option Name	Description / Valid Values
	<p data-bbox="548 369 1305 396">Specify which analyzer to use to tokenize the input string. One of these:</p> <ul data-bbox="548 422 1425 1875" style="list-style-type: none"> <li data-bbox="548 422 1425 596">• Standard—Provides a grammar-based tokenizer that contains a superset of the Whitespace and Stop Word analyzers. Understands English punctuation for breaking down words, knows words to ignore (via the Stop Word Analyzer), and performs technically case-insensitive searching by conducting lowercase comparisons. For example, the string “Precisely Software” would be returned as two tokens: “Precisely” and “Software”. <li data-bbox="548 604 1425 688">• Whitespace—Separates tokens with whitespace. Somewhat of a subset of the Standard Analyzer in that it understands word breaks in English text based on spaces and line breaks. <li data-bbox="548 697 1425 751">• StopWord—Removes articles such as "the," "and," and "a" to shrink the index size and increase performance. <li data-bbox="548 760 1425 814">• Keyword—Creates a single token from a stream of data. For example, the string “Precisely Software” would be returned as just one token “Precisely Software”. <li data-bbox="548 823 1425 907">• Russian—Supports Russian-language indexes and type-ahead services. Also supports many stop words and removes articles such as "and," "I," and "you" to shrink the index size and increase performance. <li data-bbox="548 915 1425 999">• German—Supports German-language indexes and type-ahead services. Also supports many stop words and removes articles such as "the" "and," and "a" to shrink the index size and increase performance. <li data-bbox="548 1008 1425 1092">• Danish—Supports Danish-language indexes and type-ahead services. Also supports many stop words and removes articles such as "at" "and," and "a" to shrink the index size and increase performance. <li data-bbox="548 1100 1425 1184">• Dutch—Supports Dutch-language indexes and type-ahead services. Also supports many stop words and removes articles such as "the" "and," and "a" to shrink the index size and increase performance. <li data-bbox="548 1192 1425 1276">• Finnish—Supports Finnish-language indexes and type-ahead services. Also supports many stop words and removes articles such as "is" "and," and "of" to shrink the index size and increase performance. <li data-bbox="548 1285 1425 1369">• French—Supports French-language indexes and type-ahead services. Also supports many stop words and removes articles such as "the" "and," and "a" to shrink the index size and increase performance. <li data-bbox="548 1377 1425 1461">• Hungarian—Supports Hungarian-language indexes and type-ahead services. Also supports many stop words and removes articles such as "the" "and," and "a" to shrink the index size and increase performance. <li data-bbox="548 1470 1425 1554">• Italian—Supports Italian-language indexes and type-ahead services. Also supports many stop words and removes articles such as "the" "and," and "a" to shrink the index size and increase performance. <li data-bbox="548 1562 1425 1646">• Norwegian—Supports Norwegian-language indexes and type-ahead services. Also supports many stop words and removes articles such as "the" "and," and "a" to shrink the index size and increase performance. <li data-bbox="548 1654 1425 1738">• Portuguese—Supports Portuguese-language indexes and type-ahead services. Also supports many stop words and removes articles such as "the" "and," and "a" to shrink the index size and increase performance. <li data-bbox="548 1747 1425 1831">• Spanish—Supports Spanish-language indexes and type-ahead services. Also supports many stop words and removes articles such as "the" "and," and "a" to shrink the index size and increase performance. <li data-bbox="548 1839 1425 1875">• Swedish—Supports Swedish-language indexes and type-ahead services. Also

Option Name	Description / Valid Values
	<p>supports many stop words and removes articles such as "the" "and," and "a" to shrink the index size and increase performance.</p> <ul style="list-style-type: none"> • Hindi—Supports Hindi-language indexes and type-ahead services. Also supports many stop words and removes articles such as "by" "and," and "a" to shrink the index size and increase performance.
Output Fields tab	<p>Check the Include box to select which stored fields should be included in the output.</p> <p>Note: If the input field is from an earlier stage in the dataflow and it has the same name as the store field name from the search index, the values from the input field will overwrite the values in the output field.</p>

Advanced Search Index Options

Table 15: Candidate Finder Options

Option Name	Description / Valid Values
Finder type	Select Search Index.
Name	Select the appropriate index that was created using the Write to Search Index on page 305 stage under the Advanced Matching deployed stages in Enterprise Designer.
Starting record	Enter the record number on which search results should begin. The default is 1.
Maximum results	<p>Enter the maximum number of responses you want the index search to return. The default is 10.</p> <p>Note: If the maximum results is arbitrarily large, process those in batches, using the Fetch Batch Size field.</p>

Option Name	Description / Valid Values
Fetch Batch Size	<p>If the Maximum results is arbitrarily large, enter the size of batches in which you want the results to be processed. This optimizes processing of large number of records. Default is 10000.</p> <p>The recommended Fetch Batch Size is a value lesser than Maximum results and if the Fetch Batch Size is greater than Maximum results, the records are processed in a single batch.</p> <p>Note: This field is applicable only to cluster supported search engine and not to the legacy search engine.</p>
Sort	<p>Sorts the candidate records on the basis of indexed fields while running a search query.</p> <p>Select the Sort check-box, the desired index field from the Sort by drop-down list, and select Ascending or Descending from the Order by drop-down list.</p> <p>Note: You can perform sorting only on String Fields with Keyword Analyzer and Numeric fields.</p>
Return match count	<p>Returns the total number of matches that were made. For example, if you use the default of "10" for the Maximum results field above, only 10 results will be returned. However, if you check this box, the TotalMatchCount output field will tell you how many matches were made during processing.</p>
Relevance	<p>Controls the relevance of the Index Field.</p>
Index search type	<p>Determines the type of index search you want to conduct. Select Advanced search.</p>
Add Parent button	<p>Access Parent Options.</p>
Parent options—Name	<p>Enter a name for the parent.</p>

Option Name	Description / Valid Values
Parent options—Searching method	<p>Specify how to determine if a parent is a match or a non-match. One of these:</p> <p>All true—A parent is considered a match if all children are determined to match. This method creates an "AND" connector between children.</p> <p>Any true—A parent is considered a match if at least one child is determined to match. This method creates an "OR" connector between children.</p> <p>None true—A parent is considered a match if none of the children is determined to match. This method creates a "NOT" connector between children.</p>
Add Child button	Access Child Options.
Child options—Index field	Select the index field you want to use for comparison in the advanced search.
Child options—Search type	Specifies the searching/matching criteria that determines whether the input data is searched/matched with the indexed data. All searches are case insensitive.
Child options—Input field	Select the input field you want to use for comparison in the advanced search.
Any Word/Phrase Starts With	<p>Determines whether the text contained in the search index field begins with the text that is contained in the input field.</p> <p>For example, text in the input field "tech" would be considered a match for search index fields containing "Technical", "Technology", "Technologies", "Technician" or even "National University of Technical Sciences". Likewise, a phrase in the input field "DEF Sof" would be considered a match for search index fields containing "ABC DEF Software", "DEF Software", and "DEF Software India" but it would not be a match for search index fields containing "Software DEF" or "DEF ABC Software".</p>
Contains	<p>Determines whether the search index field contains the data from the input field. This search type considers the sequence of words in the input field while searching the search index field. For example, input field data "Precisely" and "Precisely Software" would be contained in a search index field of "Precisely Software Inc."</p>
Contains All	<p>Determines whether all alphanumeric words from the input field are contained in the search index field. This search type does not consider the sequence of words in the input field while searching the search index field.</p>
Contains Any	<p>Determines whether any of the alphanumeric words from the input field is contained in the search index field.</p>

Option Name	Description / Valid Values
Contains None	Determines whether none of the alphanumeric words from the input field is contained in the search index field.
Fuzzy	<p>Determines the similarity between two alphanumeric words based on the number of deletions, insertions, or substitutions required to transform one word into another.</p> <p>Use the Maximum edits parameter to set a limit on the number of edits allowed to be considered a successful match:</p> <ul style="list-style-type: none"> • 0—Allows for no deletions, insertions, or substitutions. The input field data and the search index field data must be identical. • 1—Allows for no more than one deletion, insertion, or substitution. For example, an input field containing "Barton" will match a search index field containing "Carton". • 2—Allows for no more than two deletions, insertions, or substitutions. For example, an input field containing "Barton" will match a search index field containing "Martin". <p>The Fuzzy search type is used for single-word searches only. Click Ignore extra words to have Candidate Finder consider only the first word in the field when comparing the input field to the index field. For example, if the index field says "Xyz" and the input field says "Xyz Abc", they would not be considered a match because of "Abc". However, if you check this box, "Abc" would be ignored and with "Xyz" being the first word, the two words would be considered a match.</p>
Numeric	<p>Determines whether numbers from the input field are contained in the search index field.</p> <p>The Numeric search type is used for single-word searches only.</p> <p>Click Ignore extra words to have Candidate Finder consider only the first word in the field when comparing the input field to the index field.</p>
Pattern	<p>Determines whether the text pattern of the input field matches the text pattern of the search criteria. You can further refine the text pattern in the Pattern string field. For example, if the input field contains "nlm" and the pattern defined is "a*b?c" then it will match the following words "Neelam", "nelam", "neelum", "nilam", and so on.</p> <p>The Pattern search type is used for single-word searches only. Click Ignore extra words to have Candidate Finder consider only the first word in the field when comparing the input field to the index field.</p>

Option Name	Description / Valid Values
Proximity	<p>Determines whether words in the input fields are within a certain distance of each other.</p> <ul style="list-style-type: none"> • Define the input First input field and Second input field you want to search for in the index. • Use the Distance parameter to determine the maximum allowed distance between the words specified in the First field and Second field in order to be considered a match. <p>For example, you could successfully use this search type to look for First field "Spectrum" and Second field "Precisely" within ten words of each other in a search index field containing the sentence "Spectrum Technology Platform is a product of Precisely Software Inc."</p> <p>The Proximity search type is used for single-word searches only. Click Ignore extra words to have Candidate Finder consider only the first word in the field when comparing the input field to the index field.</p>
Range	<p>Performs an inclusive searches for terms within a range, which is specified using a Lower bound field (starting term) and an Upper bound field (ending term). All alphanumeric words are arranged lexicographically in the search index field.</p> <ul style="list-style-type: none"> • Use the Lower bound field parameter to select the field to be used as the starting term. • Use the Upper bound field parameter to select the field to be used as the ending term. <p>For example, if you searched postal codes from 20001 (defined in the Lower bound field) to 20009 (defined in the Upper bound field), the search would return all addresses with postal codes within that range.</p> <p>The Range search type is used for single-word searches only. Click Ignore extra words to have Candidate Finder consider only the first word in the field when comparing the input field to the index field.</p>
Wildcard	<p>Searches using single or multiple Wildcard characters.</p> <p>Select the Position in your input file where you are inserting the wildcard character.</p> <p>The Wildcard search type is used for single-word searches only. Click Ignore extra words to have Candidate Finder consider only the first word in the field when comparing the input field to the index field.</p>

Option Name	Description / Valid Values
Child options—Relevance factor	<p>Control the relevance of a child field by entering any positive number up to 100 here. The number can be less than "1" also; for instance, ".05" would be valid.</p> <p>The higher the boost factor, the more relevant the field will be. For example, if you want results from the Firm Name field to be more relevant than the results from other fields, select "Firm Name" from the Index field name and enter "5" here.</p> <p>Note: By default, this option is disabled. Select the check box to enable it.</p>
Ignore Blanks	<p>Clear this check-box if you want the query to take into account the blank input file fields.</p> <p>Note: By default the query ignores the blank fields.</p>
Output Fields tab	<p>Check the Include box to select which stored fields should be included in the output.</p> <p>Note: If the input field is from an earlier stage in the dataflow and it has the same name as the store field name from the search index, the values from the input field will overwrite the values in the output field.</p>

Configuring Options at Runtime

Some Candidate Finder options can be configured and passed at runtime if they are exposed as dataflow options. This enables you to run your dataflow while using different configurations. These are the available dataflow options for Candidate Finder:

- **ConnectionName**—The name of the database that contains the candidate records.
- **SearchIndexName**—The name of the search index used in the Candidate Finder dataflow.
- **StartingRecord**—The record number on which search results should begin.
- **MaximumResults**—The maximum number of responses you want the index search to return.
- **ReturnMatchCount**—The total number of matches that were made. This field is useful if you enter a lower number in the MaximumResults field but want to know the total number of matches that were made.

To define Candidate Finder options at runtime:

1. In Enterprise Designer, open a dataflow that uses the Candidate Finder stage.
2. Save and expose that dataflow.
3. Go to **EditDataflow Options**.

4. In the **Map dataflow options to stages** table, expand Candidate Finder and edit options as necessary. Check the box for the option you want to edit, then change the value in the **Default value** drop-down.
5. Optional: Change the name of the options in the **Option label** field.
6. Click **OK** twice.

Output

Table 16: Candidate Finder Outputs

Field Name	Format	Description / Valid Values
CandidateCount	String	This field indicates the total number of candidates returned during processing.
CandidateGroup	String	This field identifies a grouping of a suspect record and its candidates. Each suspect record is given a CandidateGroup number. The candidates for that suspect are given the same CandidateGroup number. For example, if John Smith is a suspect record and its candidate records are John Smith and Jon Smth, then all three records would have the same CandidateGroup value.
HasDuplicates	String	Identifies whether candidates are detected or not. The possible values are: Y - The record is suspect and has candidates N - The record is suspect and doesn't have candidates D - The record is a candidate record
TotalMatchCount	String	This field indicates the total number of matches made during the processing.
TransactionRecordType	String	One of the following: Suspect A suspect record is used as input to a query. Candidate A candidate record is a result returned from a query.

Duplicate Synchronization

Duplicate Synchronization determines which fields from a collection of records to copy to the corresponding fields of all records in the collection. You can specify the rules that records must satisfy in order to copy the field data to the other records in the collection. When processing has been completed, all records in the collection are retained.

Options

The table lists the options for the **Duplicate Synchronization** stage.

Option Name	Description / Valid Values
Group by	Specifies the field to use to create groups of records to synchronize. In cases where you have used a matching stage earlier in the dataflow, such as Interflow Match, Intraflow Match, or Transactional Match, you should select the CollectionNumber field to use the collections created by the matching stage as the groups. However, if you want to group records by some other field, choose the field here. For example, if you want to synchronize records that have the same value in the AccountNumber field, you would select AccountNumber.
Sort	If you specify a field in the Group by field, check this box to sort the records by the value in the field you chose. This option is enabled by default.

Option Name	Description / Valid Values
Advanced	<p>Click this button to specify sort performance options. By default, the sort performance options specified in Management Console, which are the default performance options for your system, are in effect. If you want to override your system's default performance options, check the Override sort performance options box then specify the values you want in these fields:</p>
In memory record limit	<p>Specifies the maximum number of data rows a sorter will hold in memory before it starts paging to disk. By default, a sort of 10,000 records or less will be done in memory and a sort of more than 10,000 records will be performed as a disk sort. The maximum limit is 100,000 records. Typically an in-memory sort is much faster than a disk sort, so this value should be set high enough so that most of the sorts will be in-memory sorts and only large sets will be written to disk.</p> <p>Note: Be careful in environments where there are jobs running concurrently because increasing the In memory record limit setting increases the likelihood of running out of memory.</p>
Maximum number of temporary files	<p>Specifies the maximum number of temporary files that may be used by a sort process. Using a larger number of temporary files can result in better performance. However, the optimal number is highly dependent on the configuration of the server running Spectrum Technology Platform. You should experiment with different settings, observing the effect on performance of using more or fewer temporary files. To calculate the approximate number of temporary files that may be needed, use this equation:</p> $\frac{(NumberOfRecords \times 2)}{InMemoryRecordLimit} = NumberOfTempFilesN$ <p>Note: The maximum number of temporary files cannot be more than 1,000.</p>
Enable compression	<p>Specifies that temporary files are compressed when they are written to disk.</p> <p>Note: The optimal sort performance settings depends on your server's hardware configuration. You can use this equation as a general guideline to produce good sort performance: $(InMemoryRecordLimit \times MaxNumberOfTempFiles \div 2) \geq TotalNumberOfRecords$</p>

Rules

Duplicate Synchronization rules determine which records should have their data copied to all other records in the collection.

To add a rule, select Rules in the rule hierarchy and click **Add Rule**

If you specify multiple rules, you will have to select a logical operator to use between each rule. Choose **And** if you want the new rule and the previous rule to both pass in order for the condition to be met. Select **Or** if you want either the previous rule or the new rule to pass in order for the condition to be met.

Option	Description
Field name	Specifies the name of the dataflow field whose value you want to evaluate to determine whether to filter the record.
Field Type	Specifies the type of data in the field. One of the following: Non-Numeric Choose this option if the field contains non-numeric data (for example, string data). Numeric Choose this option if the field contains numeric data (for example, double, float, and so on).

Option	Description
Operator	Specifies the type of comparison you want to use to evaluate the field. One of the following:
Contains	Determines if the field contains the value specified. For example, "sailboat" contains the value "boat".
Equal	Determines if the field contains the exact value specified.
Greater Than	Determines if the field value is greater than the value specified. This operation only works on numeric fields.
Greater Than Or Equal To	Determines if the field value is greater than or equal to the value specified. This operation only works on numeric fields.
Highest	Compares the field's value for all the records group and determines which record has the highest value in the field. For example, if the fields in the group contain values of 10, 20, 30, and 100, the record with the field value 100 would be selected. This operation only works on numeric fields. If multiple records are tied for the longest value, one record is selected.
Is Empty	Determines if the field contains no value.
Is Not Empty	Determines if the field contains any value.
Less Than	Determines if the field value is less than the value specified. This operation only works on numeric fields.
Less Than Or Equal To	Determines if the field value is less than or equal to the value specified. This operation only works on numeric fields.
Longest	Compares the field's value for all the records group and determines which record has the longest (in bytes) value in the field. For example, if the group contains the values "Mike" and "Michael", the record with the value "Michael" would be selected. If multiple records are tied for the longest value, one record is selected.
Lowest	Compares the field's value for all the records group and determines which record has the lowest value in the field. For example, if the fields in the group contain values of 10, 20, 30, and 100, the record with the field value 10 would be selected. This operation only works on numeric fields. If multiple records are tied for the longest value, one record is selected.
Most Common	Determines if the field value contains the value that occurs most frequently in this field among the records in the group. If two or more values are most common, no action is taken.
Not Equal	Determines if the field value is not the same as the value specified.

Option	Description
Value type	<p>Specifies the type of value you want to compare to the field's value. One of the following:</p> <p>Note: This option is not available if you select the operator Highest, Lowest, or Longest.</p> <p>Field Choose this option if you want to compare another dataflow field's value to the field.</p> <p>String Choose this option if you want to compare the field to a specific value.</p>
Value	<p>Specifies the value to compare to the field's value. If you selected Field in the Field type field, select a dataflow field. If you selected String in the Value type field, type the value you want to use in the comparison.</p> <p>Note: This option is not available if you select the operator Highest, Lowest, or Longest.</p>

Actions

Actions determine which field to copy to other records in the group. To add an action, select Actions in the Duplicate Synchronization condition tree then click the **Add Action**. Use the following options to define the action.

Option	Description
Source type	<p>Specifies the type of data to copy to other records in the group. One of the following.</p> <p>Field Choose this option if you want to copy a value from a field to the other records in the group.</p> <p>String Choose this option if you want to copy a constant value to the other records in the group.</p>
Source data	<p>Specifies the data to copy to the other records in the group. If the source type is Field, select the field whose value you want to copy to the other records in the group. If the source type is String, specify a constant value to copy to the other records in the group.</p> <p>Note: In case the source data has null value it will not be copied to the other records of the group. The other records will rather retain their original values.</p>

Option	Description
Destination	Specifies the field in the other records to which you want to copy the data specified in the Source data field. For example, if you want to copy the data to the AccountBalance field in all the other records in the group, you would specify AccountBalance.

Example of a Duplicate Synchronization Rule and Action

This Duplicate Synchronization rule and action selects the record where the match score is 100 and copies the account number AccountNumber field in all the other records in the group.

Rule

Field Name: MatchScore
 Field Type: Numeric
 Operator: Equal
 Value Type: String
 Value: 100

Action

Source Type: Field
 Source Data: AccountNumber
 Destination: NewAccountNumber

Filter

The Filter stage retains or removes records from a group of records based on the rules you specify.

Options

The following table lists the options for the Filter stage.

Option Name	Description / Valid Values
Group by	Specifies the field to use to create groups of records to filter. The Filter stage will retain one or more records from each group, depending on how you configure the stage. In cases where you have used a matching stage earlier in the dataflow, such as Interflow Match, Intraflow Match, or Transactional Match, you should select the CollectionNumber field to use the collections created by the matching stage as the groups. However, if you want to group records by some other field, choose the field here. For example, if you want to filter out all but one record from records that have the same value in the AccountNumber field, you would select AccountNumber.
Sort	If you specify a field in the Group by field, check this box to sort the records by the value in the field you chose. This option is enabled by default.

Option Name	Description / Valid Values
Advanced	<p data-bbox="548 380 1425 527">Click this button to specify sort performance options. By default, the sort performance options specified in Management Console, which are the default performance options for your system, are in effect. If you want to override your system's default performance options, check the Override sort performance options box then specify the values you want in these fields:</p> <p data-bbox="557 541 1425 793">In memory record limit Specifies the maximum number of data rows a sorter will hold in memory before it starts paging to disk. By default, a sort of 10,000 records or less will be done in memory and a sort of more than 10,000 records will be performed as a disk sort. The maximum limit is 100,000 records. Typically an in-memory sort is much faster than a disk sort, so this value should be set high enough so that most of the sorts will be in-memory sorts and only large sets will be written to disk.</p> <p data-bbox="813 810 1425 926">Note: Be careful in environments where there are jobs running concurrently because increasing the In memory record limit setting increases the likelihood of running out of memory.</p> <p data-bbox="557 968 1425 1220">Maximum number of temporary files Specifies the maximum number of temporary files that may be used by a sort process. Using a larger number of temporary files can result in better performance. However, the optimal number is highly dependent on the configuration of the server running Spectrum Technology Platform. You should experiment with different settings, observing the effect on performance of using more or fewer temporary files. To calculate the approximate number of temporary files that may be needed, use this equation:</p> $(NumberOfRecords \times 2) \div InMemoryRecordLimit = NumberOfTempFilesN$ <p data-bbox="813 1304 1425 1360">Note: The maximum number of temporary files cannot be more than 1,000.</p> <p data-bbox="557 1398 1425 1455">Enable compression Specifies that temporary files are compressed when they are written to disk.</p> <p data-bbox="634 1482 1425 1602">Note: The optimal sort performance settings depends on your server's hardware configuration. You can use this equation as a general guideline to produce good sort performance: $(InMemoryRecordLimit \times MaxNumberOfTempFiles \div 2) \geq TotalNumberOfRecords$</p>

Option Name	Description / Valid Values
Limit number of returned duplicate records	<p>Specifies the maximum number of records that are returned from each group. If you set this option to 1, you can define filter rules to determine which record in each group should be returned. If no rules are defined, the first record in each collection is returned and the rest are discarded. In this mode, the filter rules define which record will be retained.</p> <p>For example, if you define a rule where the record with the highest match score in a group is retained, and you set this option to 1, then the record with the highest match score in each group will survive and the other records in the group will be discarded.</p> <p>If you set this option to a value higher than one, you cannot specify filter rules.</p> <p>Note: In the event no records in the collection meet the defined rule criteria, then no records from the group are returned.</p>
Remove duplicates from collection	<p>Specifies to use filter rules to determine which records are removed from the collection. The remaining records in the collection are retained. When this option is selected, you must define a rule.</p> <p>Note: If a group contains only one record, the filter rules are ignored and the record is retained.</p>

Rule Options

Filter rules determine which records in a group to retain or remove. If you select the option **Limit number of returned duplicate records** then the rules determine which records survive the filter. If you select the option **Remove duplicates from collection** then the rules determine which records are removed from the dataflow.

To add a rule, select Rules in the rule hierarchy and click **Add Rule**

If you specify multiple rules, you will have to select a logical operator to use between each rule. Choose **And** if you want the new rule and the previous rule to both pass in order for the condition to be met. Select **Or** if you want either the previous rule or the new rule to pass in order for the condition to be met.

Note: You can only have one condition in a Filter stage. When you select Condition in the rule hierarchy, the buttons are grayed out.

Option	Description
Field name	Specifies the name of the dataflow field whose value you want to evaluate to determine whether to filter the record.
Field Type	Specifies the type of data in the field. One of the following: Non-Numeric Choose this option if the field contains non-numeric data (for example, string data). Numeric Choose this option if the field contains numeric data (for example, double, float, and so on).

Option	Description
Operator	Specifies the type of comparison you want to use to evaluate the field. One of the following:
Contains	Determines if the field contains the value specified. For example, "sailboat" contains the value "boat".
Equal	Determines if the field contains the exact value specified.
Greater Than	Determines if the field value is greater than the value specified. This operation only works on numeric fields.
Greater Than Or Equal To	Determines if the field value is greater than or equal to the value specified. This operation only works on numeric fields.
Highest	Compares the field's value for all the records group and determines which record has the highest value in the field. For example, if the fields in the group contain values of 10, 20, 30, and 100, the record with the field value 100 would be selected. This operation only works on numeric fields. If multiple records are tied for the longest value, one record is selected.
Is Empty	Determines if the field contains no value.
Is Not Empty	Determines if the field contains any value.
Less Than	Determines if the field value is less than the value specified. This operation only works on numeric fields.
Less Than Or Equal To	Determines if the field value is less than or equal to the value specified. This operation only works on numeric fields.
Longest	Compares the field's value for all the records group and determines which record has the longest (in bytes) value in the field. For example, if the group contains the values "Mike" and "Michael", the record with the value "Michael" would be selected. If multiple records are tied for the longest value, one record is selected.
Lowest	Compares the field's value for all the records group and determines which record has the lowest value in the field. For example, if the fields in the group contain values of 10, 20, 30, and 100, the record with the field value 10 would be selected. This operation only works on numeric fields. If multiple records are tied for the longest value, one record is selected.
Most Common	Determines if the field value contains the value that occurs most frequently in this field among the records in the group. If two or more values are most common, no action is taken.
Not Equal	Determines if the field value is not the same as the value specified.

Option	Description
Value type	<p>Specifies the type of value you want to compare to the field's value. One of the following:</p> <p>Note: This option is not available if you select the operator Highest, Lowest, or Longest.</p> <p>Field Choose this option if you want to compare another dataflow field's value to the field.</p> <p>String Choose this option if you want to compare the field to a specific value.</p>
Value	<p>Specifies the value to compare to the field's value. If you selected Field in the Field type field, select a dataflow field. If you selected String in the Value type field, type the value you want to use in the comparison.</p> <p>Note: This option is not available if you select the operator Highest, Lowest, or Longest.</p>

Example of a Filter Rule

This rule retains the record in each group with the highest value in the MatchScore field. Note that **Value** and **Value Type** options do not apply when the Operator is highest or lowest.

Field Name = MatchScore
 Field Type = Numeric
 Operator = Highest

This rule retains the record where the value in the AccountNumber is "12345".

Field Name = AccountNumber
 Field Type = Numeric
 Operator = Equals
 Value Type = String
 Value = 12345

Interflow Match

Interflow Match locates matches between similar data records across two input record streams. The first record stream is a source for suspect records and the second stream is a source for candidate records.

Using match group criteria (for example a match key), Interflow Match identifies a group of records that are potentially duplicates of a particular suspect record.

Each candidate is separately matched to the Suspect and is scored according to your match rules. If the candidate is a duplicate, it is assigned a collection number, the match record type is labeled a duplicate, and written out; unmatched unique candidates may be written out at the user's option. When Interflow Match has exhausted all candidate records in the current match group, the matched suspect record is assigned a collection number that corresponds to its duplicate record. Or, if no matches were identified, the suspect is assigned a collection number of 0 and is labeled a unique record.

Note: Interflow Match only matches suspect records to candidate records. It does not attempt to match suspect records to other suspect records as is done in Intraflow Match.

The matching process for a particular suspect may terminate before matching all possible candidates if you have set a limiter on duplicates and the limit has been exceeded for the current suspect.

The type of matching (Intraflow or Interflow) determines how express key match results translate to Candidate Match Scores. In Interflow matching, a successful Express Key match always confers a 100 MatchScore onto the Candidate. On the other hand, in Intraflow matching, the score a Candidate gains as a result of an Express Key match depends on whether the record to which that Candidate matched was a match of some other Suspect—Express Key duplicates of a Suspect will always have MatchScores of 100, whereas Express Key duplicates of another Candidate (which was a duplicate of a Suspect) will inherit the MatchScore (not necessarily 100) of that Candidate

Options

1. In the **Load match rule** field, select one of the predefined match rules which you can either use as-is or modify to suit your needs. If you want to create a new match rule without using one of the predefined match rules as a starting point, click **New**. You can only have one custom rule in a dataflow.

Note: The Dataflow Options feature in Enterprise Designer enables the match rule to be exposed for configuration at runtime.

2. Click **Group By** to select a field to use for grouping records in the match queue. Intraflow Match only attempts to match records against other records in the same match queue.
3. Select the **Sort** box to perform a pre-match sort of your input based on the field selected in the **Group By** field.
4. Click **Advanced** to specify additional sort performance options.

In memory record limit

Specifies the maximum number of data rows a sorter will hold in memory before it starts paging to disk. By default, a sort of 10,000 records or less will be done in memory and a sort of more than 10,000 records will be performed as a disk sort. The maximum limit is 100,000 records. Typically an in-memory sort is much faster than a disk sort, so this value should be set high enough so that most of the sorts will be in-memory sorts and only large sets will be written to disk.

Note: Be careful in environments where there are jobs running concurrently because increasing the **In memory record limit** setting increases the likelihood of running out of memory.

Maximum number of temporary files

Specifies the maximum number of temporary files that may be used by a sort process. Using a larger number of temporary files can result in better performance. However, the optimal number is highly dependent on the configuration of the server running Spectrum Technology Platform. You should experiment with different settings, observing the effect on performance of using more or fewer temporary files. To calculate the approximate number of temporary files that may be needed, use this equation:

$$\frac{(NumberOfRecords \times 2)}{InMemoryRecordLimit} = NumberOfTempFilesN$$

Note: The maximum number of temporary files cannot be more than 1,000.

Enable compression

Specifies that temporary files are compressed when they are written to disk.

Note: The optimal sort performance settings depends on your server's hardware configuration. You can use this equation as a general guideline to produce good sort performance: $(InMemoryRecordLimit \times MaxNumberOfTempFiles \div 2) \geq TotalNumberOfRecords$

5. Click **Express Match On** to perform an initial comparison of express key values to determine whether two records are considered a match.

Express Key matching can be a useful tool for reducing the number of compares performed and thereby improving execution speed. A loose express key results in many false positive matches. You can generate an express key as part of generating a match key through MatchKeyGenerator. See [Match Key Generator](#) on page 294 for more information.

If two records have an exact match on the express key, the candidate is considered a 100% duplicate. If two records do not match on an express key value, they are compared using the rules-based method.

To determine whether a candidate was matched using an express key, look at the value of the **ExpressKeyIdentified** field, which is either Y for a match or N for no match. Note that suspect records always have an **ExpressKeyIdentified** value of N.

6. In the **Initial Collection Number** text box, specify the starting number to assign to the collection number field for duplicate records.

The collection number identifies each duplicate record in a match queue. Unique records are assigned a collection number of 0. Each duplicate record is assigned a collection number starting with the value specified in the **Initial Collection Number** text box.

7. Select one of the following:

Option	Description
<p>Compare suspect to all candidates</p>	<p>This option matches the suspect to all candidates in the same match group (group by option) even if a duplicate is already found within the match group. For example:</p> <p>Suspect - John Smith Candidate - Bill Jones Candidate - John Smith Candidate - John Smith</p> <p>In the example, the suspect John Smith would be compared to both John smith candidates.</p> <p>Check the Return Unique Candidates box to return records within a match group from the candidate port that have been identified as unique records.</p>
<p>Stop comparing suspect against candidates after finding <i>n</i> duplicates</p>	<p>This option matches the suspect to all candidates in the same match group (group by option) but stops comparing when the user defined number of duplicates have been identified. For example, if you chose to stop comparing candidates after finding one duplicate and you had this data:</p> <p>Suspect - John Smith Candidate - Bill Jones Candidate - John Smith Candidate - John Smith</p> <p>In the example, the suspect record John Smith would stop comparing within the match group when the first John Smith candidate is identified as a duplicate.</p>

8. Click **Generate Data for Analysis** to generate match results. For more information, see [Analyzing Match Results](#) on page 111.

9. **Assign collection number 0 to unique records**, checked by default, will assign zeroes as collection numbers to unique records. Uncheck this option to generate collection numbers other than zero for unique records. The unique record collection numbers will be in sequence with any other collection numbers. For example, if your matching dataflow finds five records and the first three records are unique, the collection numbers would be assigned as shown in the first group below. If your matching dataflow finds five records and the last two are unique, the collection numbers would be assigned as shown in the second group below.

Option	Description
Collection Number	Record Type
1	Unique
2	Unique
3	Unique
4	Duplicate/Suspect
4	Duplicate/Suspect
Collection Number	Record Type
1	Duplicate/Suspect
1	Duplicate/Suspect
2	Unique
3	Unique
4	Unique

If you leave this box checked, any unique records found in your dataflow will be assigned a collection number of zero by default.

10. Select the **Return match rule name** option to include the selected match rule name in the stage output.
11. Select **Return detailed match information** if you want detailed match information to be displayed as an output for your match rule. For more information about the output fields, see [Output](#) on page 287.

Note: If you enable this field, it will hinder the overall stage performance.

12. If you are creating a new custom matching rule, see [Building a Match Rule](#) on page 69 for more information.

13. Click **Evaluate** to evaluate how a suspect record scored against candidate records. For more information, see [Interflow Match](#) on page 282.

Output

Table 17: Interflow Match Output Fields

Field Name	Description / Valid Values
CollectionNumber	Identifies a collection of duplicate records. The possible values are 1 or greater.
ExpressMatchIdentified	Indicates whether the match was obtained using the express match key. The possible values are Yes or No.
HasDuplicates	Identifies whether the record is a duplicate of another record. One of the following: Y The record is a suspect record and has duplicates. N The record is a suspect record and has no duplicates. D The record is a candidate record and is a duplicate of the suspect record. U The record is a candidate record but is not a duplicate of the suspect record.
InterflowSourceType	The possible values are input_port_0 or input_port_1
MatchRecordType	Identifies the type of match record in a collection. The possible values are: suspect The original input record that was flagged as possibly having duplicate records. duplicate A record that is a duplicate of the input record. unique A record that has no duplicates.
MatchInfo	It displays the detailed information of each match record in JSON format, for example, a match record's name, type, state as True or False, and score of the algorithm between 0-100. Note: This field is displayed only for a duplicate match record.

Field Name	Description / Valid Values
<i>MatchInfo.RuleName.IsMatch</i>	Displays the match state for the rule. Note: This field is displayed only for a duplicate match record.
<i>MatchInfo.RuleName.Score</i>	Displays the match score for the rule. The possible values are 0-100, with 0 indicating a poor match and 100 indicating a perfect match. Note: This field is displayed only for a duplicate match record.
<i>MatchInfo.RuleName.RuleNodeName.IsMatch</i>	Displays the match state for each node in the rule hierarchy, only if it has participated in the matching process, else, a blank field is displayed. Note: This field is displayed only for a duplicate match record.
<i>MatchInfo.RuleName.RuleNodeName.Score</i>	Displays the match score for each node in the rule hierarchy, only if it has participated in the matching process, else, a blank field is displayed. Note: This field is displayed only for a duplicate match record.
<i>MatchRuleName</i>	Displays the name of the match rule against which matching is performed.
<i>MatchScore</i>	Identifies the overall score between two records. The possible values are 0-100, with 0 indicating a poor match and 100 indicating an exact match.

Note: The Validate Address and Advanced Matching stages both use the MatchScore field. The MatchScore field value in the output of a dataflow is determined by the last stage to modify the value before it is sent to an output stage. If you have a dataflow that contains Validate Address and Advanced Matching stages and you want to see the MatchScore field output for each stage, use a Transformer stage to copy the MatchScore value to another field. For example, Validate Address produces an output field called MatchScore and then a Transformer stage copies the MatchScore field from Validate Address to a field called AddressMatchScore. When the matcher stage runs it populates the MatchScore field with the value from the matcher and passes through the AddressMatchScore value from Validate Address.

Intraflow Match

Intraflow Match locates matches between similar data records within a single input stream. You can create hierarchical rules based on any fields that have been defined or created in other stages of the dataflow.

Options

1. In the **Load match rule** field, select one of the predefined match rules which you can either use as-is or modify to suit your needs. If you want to create a new match rule without using one of the predefined match rules as a starting point, click **New**. You can only have one custom rule in a dataflow.

Note: The Dataflow Options feature in Enterprise Designer enables the match rule to be exposed for configuration at runtime.

2. Click **Group By** to select a field to use for grouping records in the match queue. Intraflow Match only attempts to match records against other records in the same match queue.
3. Select the **Sort** box to perform a pre-match sort of your input based on the field selected in the **Group By** field.
4. Click **Advanced** to specify additional sort performance options.

In memory record limit

Specifies the maximum number of data rows a sorter will hold in memory before it starts paging to disk. By default, a sort of 10,000 records or less will be done in memory and a sort of more than 10,000 records will be performed as a disk sort. The maximum limit is 100,000 records. Typically an in-memory sort is much faster than a disk sort, so this value should be set high enough so that most of the sorts will be in-memory sorts and only large sets will be written to disk.

Note: Be careful in environments where there are jobs running concurrently because increasing the **In memory record limit** setting increases the likelihood of running out of memory.

Maximum number of temporary files

Specifies the maximum number of temporary files that may be used by a sort process. Using a larger number of temporary files can result in better performance. However, the optimal number is highly dependent on the configuration of the server running Spectrum Technology Platform. You should experiment with different settings, observing the effect on performance of using more or fewer temporary files. To calculate the approximate number of temporary files that may be needed, use this equation:

$$\text{NumberOfRecords} \times 2 \div \text{InMemoryRecordLimit} = \text{NumberOfTempFilesN}$$

Note: The maximum number of temporary files cannot be more than 1,000.

Enable compression Specifies that temporary files are compressed when they are written to disk.

Note: The optimal sort performance settings depends on your server's hardware configuration. You can use this equation as a general guideline to produce good sort performance: $(InMemoryRecordLimit \times MaxNumberOfTempFiles \div 2) \geq TotalNumberOfRecords$

5. Click **Express Match On** to perform an initial comparison of express key values to determine whether two records are considered a match.

You can generate an express key as part of generating a match key through MatchKeyGenerator. See [Match Key Generator](#) on page 294 for more information.

6. In the **Initial Collection Number** text box, specify the starting number to assign to the collection number field for duplicate records.

The collection number identifies each duplicate record in a match queue. Unique records are assigned a collection number of 0. Each duplicate record is assigned a collection number starting with the value specified in the **Initial Collection Number** text box.

7. Click **Sliding Window** to enable this matching method. For more information about Sliding Window, see [Sliding Window Matching Method](#) on page 292
8. Click **Generate Data for Analysis** to generate match results. For more information, see [Analyzing Match Results](#) on page 111.

9. **Assign collection number 0 to unique records**, checked by default, will assign zeroes as collection numbers to unique records. Uncheck this option to generate collection numbers other than zero for unique records. The unique record collection numbers will be in sequence with any other collection numbers. For example, if your matching dataflow finds five records and the first three records are unique, the collection numbers would be assigned as shown in the first group below. If your matching dataflow finds five records and the last two are unique, the collection numbers would be assigned as shown in the second group below.

Option	Description
Collection Number	Record Type
1	Unique
2	Unique
3	Unique

Option	Description
4	Duplicate/Suspect
4	Duplicate/Suspect
Collection Number	Record Type
1	Duplicate/Suspect
1	Duplicate/Suspect
2	Unique
3	Unique
4	Unique

If you leave this box checked, any unique records found in your dataflow will be assigned a collection number of zero by default.

10. Select the **Return match rule name** option to include the selected match rule name in the stage output.
11. Select **Return detailed match information** if you want detailed match information to be displayed as an output for your match rule. For more information about the output fields, see **Output** on page 292.

Note: If you enable this field, it will hinder the overall stage performance.

12. For information about modifying the other options, see **Building a Match Rule** on page 69.
13. Click **Evaluate** to evaluate how a suspect record scored against candidate records. For more information, see **Interflow Match** on page 282.

Default Matching Method

Using group by (match group) set by the user, the matcher identifies groups of records that might potentially be duplicates of one another. The matcher then proceeds through each record in the group; if the record matches an existing Suspect, the record is considered a Duplicate of that suspect, assigned a Score, CollectionNumber, and MatchRecordType (Duplicate), and eliminated from the match. If, on the other hand, the record matches no existing Suspect within the match group, the record becomes a new Suspect, in that it is added to the current Match group so that it can be

matched against by subsequent records. When the matcher has exhausted all records in the current Match group, it eliminates all Suspects from the match, labeling the Match Record type as Unique and assigning a collection number of 0. Those Suspects with a least one duplicate will retain a Match Record Type of Suspect and is assigned the same collection number as its matched duplicate record. Finally, when all records within a match group have been written to the output. A new match group is compared.

Note: The Default Matching Method will only compare records that are within the same match group.

The type of matching (Intraflow or Interflow) determines how express key match results translate to Candidate Match Scores. In Interflow matching, a successful Express Key match always confers a 100 MatchScore onto the Candidate. On the other hand, in Intraflow matching, the score a Candidate gains as a result of an Express Key match depends on whether the record to which that Candidate matched was a match of some other Suspect—Express Key duplicates of a Suspect will always have MatchScores of 100, whereas Express Key duplicates of another Candidate (which was a duplicate of a Suspect) will inherit the MatchScore (not necessarily 100) of that Candidate

Sliding Window Matching Method

The sliding window algorithm is an algorithm which sequentially fills a pre determined buffer size called a window with the corresponding amount of data rows. As each row is added to the window it's compared to each item already contained in the window. If a match with an item is determined then both the driver record (the new item to add to the window) and the candidates (items already in the window) is given the same group ID. This comparison is continued until the driver record has been compared to all items contained within the window.

As new drivers are added the window will eventually reach its predetermined capacity. At this point the window will slide, hence the term Sliding Window. Sliding simply means that the window buffer will remove and write the oldest item in the window as it adds the newest driver record to the window.

Output

Table 18: Intraflow Match Output

Field Name	Description / Valid Values
CollectionNumber	Identifies a collection of duplicate records. The possible values are 1 or greater.
ExpressMatchIdentified	Indicates whether the match was obtained using the express match key. The possible values are Yes or No.

Field Name	Description / Valid Values
MatchRecordType	<p>Identifies the type of match record in a collection. The possible values are:</p> <p>Suspect A record that other records are compared to in order to determine if they are duplicates of each other. Each collection has one and only one suspect record.</p> <p>Duplicate A record that is a duplicate of the suspect record.</p> <p>Unique A record that has no duplicates.</p>
MatchInfo	<p>It displays the detailed information of each match record in JSON format, for example, a match record's name, type, state as True or False, and score of the algorithm between 0-100.</p> <p>Note: This field is displayed only for a duplicate match record.</p>
MatchInfo.RuleName.IsMatch	<p>Displays the match state for the rule.</p> <p>Note: This field is displayed only for a duplicate match record.</p>
MatchInfo.RuleName.Score	<p>Displays the match score for the rule. The possible values are 0-100, with 0 indicating a poor match and 100 indicating a perfect match.</p> <p>Note: This field is displayed only for a duplicate match record.</p>
MatchInfo.RuleName. RuleNodeName.IsMatch	<p>Displays the match state for each node in the rule hierarchy, only if it has participated in the matching process, else, a blank field is displayed.</p> <p>Note: This field is displayed only for a duplicate match record.</p>
MatchInfo.RuleName. RuleNodeName.Score	<p>Displays the match score for each node in the rule hierarchy, only if it has participated in the matching process, else, a blank field is displayed.</p> <p>Note: This field is displayed only for a duplicate match record.</p>
MatchRuleName	<p>Displays the name of the match rule against which matching is performed.</p>
MatchScore	<p>Identifies the overall score between two records. The possible values are 0-100, with 0 indicating a poor match and 100 indicating an exact match.</p>

Note: The Validate Address and Advanced Matching stages both use the MatchScore field. The MatchScore field value in the output of a dataflow is determined by the last stage to modify the value before it is sent to an output stage. If you have a dataflow that contains Validate Address and Advanced Matching stages and you want to see the MatchScore field output for each stage, use a Transformer stage to copy the MatchScore value to another field. For example, Validate Address produces an output field called MatchScore and then a Transformer stage copies the MatchScore field from Validate Address to a field called AddressMatchScore. When the matcher stage runs it populates the MatchScore field with the value from the matcher and passes through the AddressMatchScore value from Validate Address.

Match Key Generator

Match Key Generator creates a non-unique key for each record, which can then be used by matching stages to identify groups of potentially duplicate records. Match keys facilitate the matching process by allowing you to group records by match key and then only comparing records within these groups.

The match key is created using rules you define and is comprised of input fields. Each input field specified has a selected algorithm that is performed on it. The result of each algorithm is then concatenated to create a single match key field.

In addition to creating match keys, you can also create express match keys to be used later in the dataflow by an Intraflow Match stage or an Interflow Match stage.

You can create multiple match keys and express match keys.

For example, if the incoming record is:

First Name - Fred
 Last Name - Mertz
 Postal Code - 21114-1687
 Gender Code - M

And you define a match key rule that generates a match key by combining data from the record like this:

Input Field	Start Position	Length
Postal Code	1	5
Postal Code	7	4
Last Name	1	5
First Name	1	5

Input Field	Start Position	Length
Gender Code	1	1

Then the key would be:

211141687MertzFredM

Input

The input is any field in the source data.

Options

To define Match Key Generator options click the **Add** button. The **Match Key Field** dialog displays.

Note: The Dataflow Options feature in Enterprise Designer enables Match Key Generator to be exposed for configuration at runtime.

Table 19: Match Key Generator Options

Option Name	Description and Valid Values
Algorithm	

Option Name Description and Valid Values

Specifies one of these algorithms to use to generate the match key:

- Consonant** Returns specified fields with consonants removed.
- Double Metaphone** Returns a code based on a phonetic representation of their characters. Double Metaphone is an improved version of the Metaphone algorithm, and attempts to account for the many irregularities found in different languages.
- Koeln** Indexes names by sound as they are pronounced in German. Allows names with the same pronunciation to be encoded to the same representation so that they can be matched, despite minor differences in spelling. The result is always a sequence of numbers; special characters and white spaces are ignored. This option was developed to respond to limitations of Soundex.
- MD5** A message digest algorithm that produces a 128-bit hash value. This algorithm is commonly used to check data integrity.
- Metaphone** Returns a Metaphone coded key of selected fields. Metaphone is an algorithm for coding words using their English pronunciation.
- SpanishMetaphone** Returns a Metaphone coded key of selected fields for the Spanish language. This metaphone algorithm codes words using their Spanish pronunciation.
- Metaphone 3** Improves upon the Metaphone and Double Metaphone algorithms with more exact consonant and internal vowel settings that allow you to produce words or names more or less closely matched to search terms on a phonetic basis. Metaphone 3 increases the accuracy of phonetic encoding to 98%. This option was developed to respond to limitations of Soundex.
- Nysiis** Phonetic code algorithm that matches an approximate pronunciation to an exact spelling and indexes words that are pronounced similarly. Part of the New York State Identification and Intelligence System. Say, for example, that you are looking for someone's information in a database of people. You believe that the person's name sounds like "John Smith", but it is in fact spelled "Jon Smyth". If you conducted a search looking for an exact match for "John Smith" no results would be returned. However, if you index the database using the NYSIIS algorithm and search using the NYSIIS algorithm again, the correct match will be returned because both "John Smith" and "Jon Smyth" are indexed as "JAN SNATH" by the algorithm.
- Phonix** Preprocesses name strings by applying more than 100 transformation rules to single characters or to sequences of several characters. 19 of those rules are applied only if the characters are at the beginning of the string, while 12 of the rules are applied only if they are at the middle of the string, and 28 of the rules are applied only if they are at the end of the string. The

Option Name	Description and Valid Values
	<p>transformed name string is encoded into a code that is comprised by a starting letter followed by three digits (removing zeros and duplicate numbers). This option was developed to respond to limitations of Soundex; it is more complex and therefore slower than Soundex.</p>
	<p>Sonnex This algorithm determines the similarity between two French-language strings based on the phonetic representation of their characters. It returns a Sonnex coded key of the selected fields.</p>
	<p>Soundex Returns a Soundex code of selected fields. Soundex produces a fixed-length code based on the English pronunciation of a word.</p>
	<p>Substring Returns a specified portion of the selected field.</p>
Field name	<p>Specifies the field to which you want to apply the selected algorithm to generate the match key. For example, if you select a field called LastName and you choose the Soundex algorithm, the Soundex algorithm would be applied to the data in the LastName field to produce a match key.</p>
Start position	<p>Specifies the starting position within the specified field. Not all algorithms allow you to specify a start position.</p>
Length	<p>Specifies the length of characters to include from the starting position. Not all algorithms allow you to specify a length.</p>
Remove noise characters	<p>Removes all non-numeric and non-alpha characters such as hyphens, white space, and other special characters from an input field.</p>
Sort input	<p>Sorts all characters in an input field or all terms in an input field in alphabetical order.</p> <p>Characters Sorts the characters values from an input field prior to creating a unique ID.</p> <p>Terms Sorts each term value from an input field prior to creating a unique ID.</p>

If you add multiple match key generation algorithms, you can use the **Move Up** and **Move Down** buttons to change the order in which the algorithms are applied.

Generating an Express Match Key

Enable the **Generate Express Match Key** option and click **Add** to define an express match key to be used later in the dataflow by an Intraflow Match stage or an Interflow Match stage.

If the **Generate Express Match Key** option is enabled and the **Express match key on** option is selected in a downstream Interflow Match stage or Intraflow Match stage, the match attempt is first made using the express match key created here. If two records' express match keys match, then the record is considered a match and no further processing is attempted. If the records' express match keys do not match, then the match rules defined in Interflow Match or Intraflow Match are used to determine if the records match.

Output

Table 20: Match Key Generator Output

Field Name	Description / Valid Values
ExpressMatchKey	A value indicating the match level. If the express match key is a match, the score is 100. If the express match key does not match, then a score of 0 is returned.
MatchKey	The key generated to identify records.

Private Match

Private Match enables two entities to compare datasets and identify common records without compromising sensitive information. For example, two companies could be interested in launching a joint marketing campaign. Each company has its own database containing customer information, and the companies want to determine which customers shop at both companies to use a more targeted approach in the campaign. However, to ensure data security and comply with privacy regulations, the companies do not wish to share these databases with each other or to give them to a third party to conduct a match. The private match feature makes it possible for the two databases to be matched against each other without breaching security or breaking privacy laws.

Private Match is used in one of three modes:

- **Encrypt mode**—The first user inputs his data, and an index field and match field are extracted and encrypted. A public key and a displacement table containing the first user's data are generated for the second user, and a private key is generated for the first user to use later.

- Private Match mode—The second user inputs his data and the first user's encrypted data, provides the public key and displacement table, and performs a match. A file containing the matched data is generated to be sent to the first user.
- Decrypt mode—The first user inputs the second user's encrypted data, provides the private key, and generates output containing a matched index of both user's data.

By using the encrypt function (Encrypt mode) the security is retained while a match function is performed (Private Match mode), and then a decrypt function shows the output of the matched data (Decrypt mode). All files generated and shared between users are encrypted and unreadable.

Input

Input requirements for the Private Match stage vary depending on the task you are performing:

- Encrypt mode—A file containing the first user's data must be attached to the input port. This can be a text file, database, or almost any kind of source file.
- Private Match mode—A file containing the second user's data must be attached to the first input port; this can also be a text file, database, or almost any kind of source file. The encrypted data generated by the first user must be attached to the second input port.
- Decrypt mode—The output file generated by the second user.

Options

Options for the Private Match stage vary depending on the task you are performing.

Note: If you upgrade to version 11.0 or later from version 10.x and you produced private keys, you will need to regenerate those keys because keys generated in version 10.x are not compatible with versions 11.0 and later.

Encrypt Mode

1. Select the **Encrypt** operation.
2. Select the **index field** that provides a unique ID for each record in the file. The unique ID must be a numeric value.
3. Select the **match field** that should be used to match against the second user's data.
4. Specify the path to and name of the **Public key** file that will be created when you run the job.
5. Specify the path to and name of the **Key Store** file that will be created when you run the job.
6. Specify the path to and name of the **Displacement table** file that will be created when you run the job.
7. Enter a name for the **output column** that will contain the encrypted data in the output file that is sent to the second user.
8. Press **OK**.

Private Match Mode

1. Select the **Private Match** operation.
2. Select the **index field** that provides a unique ID for each record in the file. The unique ID must be a numeric value.
3. Select the **match field** that should be used to match against the first user's data.
4. Enter the name of the **encrypted data field** that the first user entered in step 7 of the Encrypt Mode instructions.
5. Navigate to the **Public key** file.
6. Navigate to the **Displacement table** file.
7. Enter a name for the **output column** that will contain the encrypted data in the output file that is sent to the first user.

Decrypt Mode

1. Select the **Decrypt** operation.
2. Select the **encrypted data field** entered in step 7 of the Private Match Mode instructions.
3. Navigate to the **Key Store** file.
4. Select the **output column** that will contain the decrypted data in the output file. The format of the data in this field is the matched index of the first user's data and the matched index of the second users' data, separated by a comma character (,), as in the following:

```
User1Data,User2Data
```

Output

Output requirements for the Private Match stage vary depending on the task you are performing:

- Encrypt mode—A file generated by the Write to File stage that contains the first user's encrypted data.
- Private Match mode—A file generated by the Write to File stage that contains encrypted information about the match results.
- Decrypt mode—A file generated by the Write to File stage that contains the matched index of both users' data.

Transactional Match

Transactional Match matches suspect records against candidate records that are returned from the Candidate Finder stage. Transactional Match uses matching rules to compare the suspect record to all candidate records with the same candidate group number (assigned in Candidate Finder) to identify duplicates. If the candidate record is a duplicate, it is assigned a collection number, the match

record type is labeled a Duplicate, and the record is then written out. Any unmatched candidates in the group are assigned a collection number of 0, labeled as Unique and then written out as well.

Note: Transactional Match only matches suspect records to candidates. It does not attempt to match suspect records to other suspect records as is done in Intraflow Match.

Options

1. In the **Load match rule** field, select one of the predefined match rules which you can either use as-is or modify to suit your needs. If you want to create a new match rule without using one of the predefined match rules as a starting point, click **New**. You can only have one custom rule in a dataflow.

Note: The Dataflow Options feature in Enterprise Designer enables the match rule to be exposed for configuration at runtime.

2. Select **Return unique candidates** if you want unique candidate records to be included in the output from the stage.
3. Select **Generate data for analysis** if you want to use the Match Analysis tool to analyze the results of the dataflow. For more information, see [Analyzing Match Results](#) on page 111.
4. Select **Return detailed match information** if you want detailed match information to be displayed as an output for your match rule.

Note: For detailed information about the output fields, see [Output](#) on page 303.

5. For information about modifying the other options, see [Building a Match Rule](#) on page 69.
6. Click **Evaluate** to evaluate how a suspect record scored against candidate records. For more information, see [Interflow Match](#) on page 282.

Output

Table 21: Transactional Match Output

Field Name	Description / Valid Values
HasDuplicates	<p>Identifies whether the record is a duplicate of another record. One of the following:</p> <p>Y The record is a suspect record and has duplicates.</p> <p>N The record is a suspect record and has no duplicates.</p> <p>D The record is a candidate record and is a duplicate of the suspect record.</p> <p>U The record is a candidate record but is not a duplicate of the suspect record.</p>
DuplicateCount	It displays the number of duplicates within a match group.
MatchGroupSize	It determines the size and displays the number of records in a match group.
MatchRecordType	<p>Identifies the type of match record in a collection. The possible values are:</p> <p>Suspect The original input record that was flagged as possibly having duplicate records.</p> <p>Duplicate A record that is a duplicate of the input record.</p> <p>Unique A record that has no duplicates.</p>
MatchInfo	It displays the detailed information of each match record in JSON format, for example, a match record's name, type, state as True or False, and score between 0-100.
MatchScore	Identifies the overall score between two records. The possible values are 0-100, with 0 indicating a poor match and 100 indicating an exact match.

Field Name	Description / Valid Values
<code>MatchInfo.RuleName.IsMatch</code>	<p>Displays the match state for the rule as <code>True</code> or <code>False</code>, with <code>True</code> representing a match and <code>False</code> indicating no matches. This information is displayed if you select the Return detailed match information option in the Transactional Match Options screen.</p> <p>Note: The value in this field is derived from the match state of the child nodes based on the defined parent rule configuration.</p>
<code>MatchInfo.RuleName.Score</code>	<p>Displays the match score for the rule. The possible values are 0-100, with 0 indicating a poor match and 100 indicating a perfect match. This information is displayed if you select the Return detailed match information option in the Transactional Match Options screen.</p> <p>Note: The score in this field is derived from the scores of the child nodes on the basis of the defined parent rule configuration.</p>
<code>MatchInfo.RuleName.RuleNodeName.IsMatch</code>	<p>Displays the match state for each node in the rule hierarchy as <code>True</code> or <code>False</code>, with <code>True</code> representing a match and <code>False</code> indicating no matches. <code>RuleNodeName</code> is a variable, which is replaced by the hierarchical node names in your match rules.</p> <p>Each node in the rule hierarchy outputs this field, if you have selected the Return detailed match information option in the Transactional Match Options screen.</p>
<code>MatchInfo.RuleName.RuleNodeName.Score</code>	<p>Displays the match score for each node in the rule hierarchy. <code>RuleNodeName</code> is a variable, which is replaced by the hierarchical node names in your match rules.</p> <p>Each node in the rule hierarchy outputs this field, if you have selected the Return detailed match information option in the Transactional Match Options screen.</p> <p>The possible values are 0-100, with 0 indicating a poor match and 100 indicating a perfect match.</p>

Note: The Validate Address and Advanced Matching stages both use the MatchScore field. The MatchScore field value in the output of a dataflow is determined by the last stage to modify the value before it is sent to an output stage. If you have a dataflow that contains Validate Address and Advanced Matching stages and you want to see the MatchScore field output for each stage, use a Transformer stage to copy the MatchScore value to another field. For example, Validate Address produces an output field called MatchScore and then a Transformer

stage copies the MatchScore field from Validate Address to a field called AddressMatchScore. When the matcher stage runs it populates the MatchScore field with the value from the matcher and passes through the AddressMatchScore value from Validate Address.

Write to Search Index

Write to Search Index enables you to create a full-text index based on the data coming in to the stage. Having this data in a dedicated search index results in quicker response time when you conduct searches against the index from other Spectrum Technology Platform stages. Full-text-search indexes are preferable to relational databases when you have a great deal of free-form text data that needs to be searched or categorized or if you support a high volume of interactive, text-based queries.

Write to Search Index uses an analyzer to break input text into small indexing elements called tokens. It then extracts search index terms from those tokens. The type of analyzer used—the manner in which input text is broken into tokens—determines how you will then be able to search for that text. For example, the *Keyword* analyzer always does an exact match and tokenizes the whole string as single token, while the *Standard* analyzer breaks the string to create tokens. The *Stop Word* analyzer is all the more sophisticated and removes articles, such as "a" or "the" from the string before tokenizing, thus shrinking the index size.

Search indexes support the near real time feature, allowing indexes to be updated almost immediately, without the need to close and rebuild the stage using the search index.

For information about the search index tasks you can perform through the command line, see Search Indexes in [Administration Utility](#) section of the Administration Guide.

Options

1. In Enterprise Designer, double-click the Write to Search Index stage on the canvas.
2. Enter a **Name** for the index.
3. Select a **Write mode**. When you regenerate an index, you have options related to how the new data should affect the existing data.
 - **Create or Overwrite**—New data will overwrite the existing data and the existing data will no longer be in the index.
 - **Update or Append**—New data will overwrite existing data, and any new data that did not previously exist will be added to the index.
 - **Append**—New data will be added to the existing data and the existing data will remain in tact.
 - **Delete**—Data for the selected field will be deleted from the search index.
4. Select the **Key field** on the basis of which you want to **Update or Append** or **Delete** the records.
 - In case of **Create or Overwrite** mode, the **Key field** needs to be unique for Elastic search indexes (used in a distributed environment). If you leave the field blank, all the records get stored in the index irrespective of any duplication. However, you will not be able to perform

any write operation, such as update, append, and delete on this index. The following table explains the indexing behavior if the **Key field** is non-unique for Lucene and Elastic search indexes.

Write mode	Key field	Lucene search index	Elastic search index
Create or Overwrite	Duplicate records with same Key field	All the records are stored. Note: The duplicate records with same key field get overwritten as soon as you run the update operation.	All duplicate records with the same Key field are overwritten .
Update or Append	Duplicate records with same Key field	Duplicates are overwritten.	Duplicates are overwritten.

5. Check the **Batch commit** box if you want to specify the number of records to commit in a batch while creating the search index. Then enter that number in the **Batch size** field. Default is 5000.
6. Select an **Analyzer** to build:
 - **Standard**—Provides a grammar-based tokenizer that contains a superset of the Whitespace and Stop Word analyzers. Understands English punctuation for breaking down words, knows words to ignore (via the Stop Word Analyzer), and performs technically case-insensitive searching by conducting lowercase comparisons. For example, the string “Precisely Software” would be returned as two tokens: “precisely” and “software”. For a comparison of Standard and Keyword analyzers, see [Standard and Keyword Analyzer](#) on page 309.
 - **Whitespace**—Separates tokens with whitespace. Somewhat of a subset of the Standard Analyzer in that it understands word breaks in English text based on spaces and line breaks.
 - **Stop Word**—Removes articles such as "the," "and," and "a" to shrink the index size and increase performance.
 - **Keyword**—Creates a single token from a stream of data and keeps it as is. For example, the string “Precisely Software” would be returned as just one token “Precisely Software”. For a comparison of Standard and Keyword analyzers, see [Standard and Keyword Analyzer](#) on page 309.
 - **Russian**—Supports Russian-language indexes and type-ahead services. Also supports many stop words and removes articles such as "and," "I," and "you" to shrink the index size and increase performance.
 - **German**—Supports German-language indexes and type-ahead services. Also supports many stop words and removes articles such as "the" "and," and "a" to shrink the index size and increase performance.

- Danish—Supports Danish-language indexes and type-ahead services. Also supports many stop words and removes articles such as "at" "and," and "a" to shrink the index size and increase performance.
 - Dutch—Supports Dutch-language indexes and type-ahead services. Also supports many stop words and removes articles such as "the" "and," and "a" to shrink the index size and increase performance.
 - Finnish—Supports Finnish-language indexes and type-ahead services. Also supports many stop words and removes articles such as "is" "and," and "of" to shrink the index size and increase performance.
 - French—Supports French-language indexes and type-ahead services. Also supports many stop words and removes articles such as "the" "and," and "a" to shrink the index size and increase performance.
 - Hungarian—Supports Hungarian-language indexes and type-ahead services. Also supports many stop words and removes articles such as "the" "and," and "a" to shrink the index size and increase performance.
 - Italian—Supports Italian-language indexes and type-ahead services. Also supports many stop words and removes articles such as "the" "and," and "a" to shrink the index size and increase performance.
 - Norwegian—Supports Norwegian-language indexes and type-ahead services. Also supports many stop words and removes articles such as "the" "and," and "a" to shrink the index size and increase performance.
 - Portuguese—Supports Portuguese-language indexes and type-ahead services. Also supports many stop words and removes articles such as "the" "and," and "a" to shrink the index size and increase performance.
 - Spanish—Supports Spanish-language indexes and type-ahead services. Also supports many stop words and removes articles such as "the" "and," and "a" to shrink the index size and increase performance.
 - Swedish—Supports Swedish-language indexes and type-ahead services. Also supports many stop words and removes articles such as "the" "and," and "a" to shrink the index size and increase performance.
 - Hindi—Supports Hindi-language indexes and type-ahead services. Also supports many stop words and removes articles such as "by" "and," and "a" to shrink the index size and increase performance.
7. To update the analyzer of all fields present in the list, Select an analyzer from **update Analyzers to...** drop down.
 8. To reload the schema from the server, click **Reload Schema**.

Note: You can change the field name by typing the new name directly in the **Fields** column. However, you cannot change the **Stage Fields** name or the **Type**.
 9. To selectively add/remove fields from your input source, click **Quick Add**. The **Quick Add** pop-up window displays a list of all the fields from the input source. Select the fields that you want to add and click **OK**.

10. Select the field(s) whose data you want to store. For example, using an input file of addresses, you could index just the Postal Code field but choose to store the remaining fields (such as Address Line 1, City, State) so the entire address is returned when a match is found using the index search.
11. Select the field(s) whose data you want to be added to the index for a search query.

Note: If you want to delete certain fields, select those and click **Delete**.

12. If necessary, change the analyzer for any field that should use something other than what you selected in the Analyzer field.
13. Click **OK**.

Output

Write to Search Index has one output port, called **Error Port**, which collects data for records that could not be processed and added to the search index. Records that pass through the error port into the sink are considered malformed.

Records are forwarded to the error port if the value of the **KeyField** in the input record is blank or null. For the search index engines that support standalone searches (the legacy index engine), the error port functions only for the Update and Delete operations. However, for distributed support index engine, it collects malformed records for all the four operations: Create, Update, Delete, and Append.

Search Index Management

The **Search Index Management** tool enables you to perform the following tasks:

- Add/remove fields in an existing search index
- Delete one or more search indexes

To add/remove fields to an existing search index:

1. Select **Tools > Search Index Management**. The **Search Index Management** pop-up window is displayed showing the existing search indexes.
2. Select the index to be modified, and click **Modify**. The **Modify Index** page is displayed showing a list of all the fields in the search index with their details.
3. To delete a field, select it from the list, and click **Remove field**.
4. To add a new field to the list, click **Add field**. The **Add Field** pop-up window is displayed, which allows you to add the following details:
 - **Field name:** Enter the name of the field.
 - **Type:** Select the field type. The options are: string, integer, long, float, and double.
 - **Index:** Select the check-box to add it to the index for search query.
 - **Store:** Select the check-box to store the field.
 - **Analyzer:** From the list of options, select the analyzer you want to use for indexing the field.

Note: This option will be enabled only if you have selected to index this field.

To delete an existing search index:

1. Select **Tools > Search Index Management**.
2. Select the search indexes you want to delete.
3. Click **Delete**.
4. Click **Close**.

Note: You can also delete a search index by using the Administration Utility. The command is `index delete --d IndexName`, where *IndexName* is the name of the index you want to delete.

Standard and Keyword Analyzer

Before choosing between the Standard and Keyword analyzers you should be aware of the following difference in the behavior of these two analyzers:

- **Standard analyzer:** It breaks the string to tokenize it and also converts all its tokens to lower case.
- **Keyword analyzer:** It tokenizes the whole string and keeps it as is (does not change the case).

Example:

Let us take an example of *contains any* algorithm and assume the input is P O. In this case, the **Standard** analyzer will return both P O and P records, while the **Keyword** analyzer will return only the P O records (shown below).

Search result with Standard analyzer:

Point		Point 2					
SearchTerm	CandidateGroup	HasDuplicates	Index	IndexField	Term	TransactionRecordType	
PO	1	Y				Suspect	
PO	1	D	PO	4	Test4	Candidate	
PO	1	D	P	1	Test 1	Candidate	

Search result with Keyword analyzer:

Point		Point 2					
SearchTerm	CandidateGroup	HasDuplicates	Index	IndexField	Term	TransactionRecordType	
PO	1	Y				Suspect	
PO	1	D	PO	4	Test4	Candidate	

Data Stewardship

Introduction

Data Stewardship stages identify and manage exception records that Spectrum Technology Platform could not confidently process to allow manual review by a data steward.

Data Stewardship provides three stages in Enterprise Designer.

- **Exception Monitor**—A stage that evaluates records against a set of conditions to determine if the record requires manual review by a data steward. When records meet those conditions, this stage can send an email notifying recipients of the exceptions. An approval flow can be added to a condition to require acceptance by reviewers in addition to the data steward.
- **Write Exceptions**—A stage that writes the exception records to the exception repository. Once exception records are in the exception repository they are available for review by a data steward and (if there is an approval flow) acceptance by other reviewers.
- **Read Exceptions**—A stage that reads records from the exception repository into a dataflow. This stage allows you to reprocess exception records that have been corrected by a data steward.

Exception Monitor

The Exception Monitor stage evaluates records against a set of conditions to determine if the record requires manual review by a data steward. Exception Monitor enables you to route records that Spectrum Technology Platform could not successfully process to a manual review tool (the Data Stewardship Portal).

In the stage settings, you can set conditions that determine if records require manual review. You can configure notifications to email addresses when those conditions have been met a certain number of times. You can assign an approval flow type to a condition to create an approval flow.

For more information on exception processing, see [Data Stewardship Portal](#).

Input

Exception Monitor takes any record as input. If the input data does not contain a field called "CollectionNumber" the **Return all records in exception's group** option will be disabled.

Note: Exception Monitor cannot monitor fields that contain complex data such as lists or geometry objects.

Output

Exception Monitor returns records in two ports. The success port transmits records that do not meet any of the conditions defined in the Exception Monitor stage. The exception port transmits records that match one or more exception conditions. The exception port may also include non-exception records if you enable the option **Return all records in exception's group**. Exception Monitor only evaluates conditions. It does not modify fields within a record.

Reference

Conditions tab

This tab is displayed in the **Exception Monitor Options** dialog box.

Stop evaluating when a condition is met Specifies whether to continue evaluating a record against the remaining conditions once a condition is met. Enabling this option may improve performance because it potentially reduces the number of evaluations that the system has to perform. However, if not all conditions are evaluated you will lose some degree of completeness in the exception reports shown in the Data Stewardship Portal. For example, if you define three conditions (Address Completeness, Name Confidence, and Geocode Confidence) and a record meets the criteria defined in Address Completeness, and you enable this option, the record would not be evaluated against Name Confidence and Geocode Confidence. If the record also qualifies as an exception because it matches the Name Confidence condition, this information would not be captured. Instead the record would be reported as having only an Address Completeness problem, instead of both an Address Completeness and Name Confidence problem.

Conditions list The list box on this tab shows all conditions that have been defined for the stage in tabular view. The conditions are evaluated in the order that they appear in the list. The list displays the following information for each condition.

- **Name**—The name for the condition. This is typically a meaningful name created by the user who created the condition.
- **Approval Flow Type**—The name of an approval flow type defined on the Management Console **Data Stewardship Settings** page.
- **Domain**—Specifies the kind of data being evaluated by the condition. This is used solely for reporting purposes to show which types of exceptions occur in the data.
- **Metric**—Specifies the metric that this condition measures. This is used solely for reporting purposes to show which types of exceptions occur in your data.
- **Assign to**—The user to whom the exception records meeting this condition are assigned. This user is referred to as the *data steward*. If this setting is empty, exception records from a job are automatically assigned to the user who ran the job.

Add Click this button to define a new condition.

Modify Click this button to edit the currently selected condition in the list.

- Remove** Click the button to remove the currently selected condition from the list.
- Move Up/Move Down** Click these buttons to order conditions that appear in the list. Conditions are evaluated in the order that they are displayed in the table. You can use these buttons to arrange selections when **Stop evaluating when a condition is met** is checked to assure that certain conditions more likely to be evaluated.

Configuration tab

This tab is displayed in the **Exception Monitor Options** dialog box.

- Disable exception monitor** Turns Exception Monitor on or off. If you disable Exception Monitor, records will simply pass through the stage and no action will be taken. This is similar in effect to removing Exception Monitor from the dataflow.
- Stop job after reaching exception limit** Specifies whether to halt job execution when the specified number of records meet the exception conditions.
- Maximum number of exception records** If **Stop job after reaching exception limit** is selected, use this field to specify the maximum number of exception records to allow before halting job execution. For example, if you specify 100, the job will stop once the 101st exception record is encountered.
- Report only (do not create exceptions)** Enables you to track records that meet exception conditions and reports those statistics on the Data Quality Performance page in the Data Stewardship Portal, but does not create exceptions for those records.
- Return all records in exception's group** Specifies whether to return all records belonging to an exception record's group instead of just the exception record. For example, a match group (based on a MatchKey) contains four records. One is the Suspect record, one is a duplicate that scored 90, and two are unique records that scored 80 and 83. If you have a condition that says that any record with a MatchScore between 80 and 89 is an exception, by default just the records with a match score of 80 and 83 would be sent to the exception port. However, if you enable this option, all four records would be sent to the exception port.
- Enable this option if you want data stewards to be able to compare the exception record to the other records in the group. By comparing all the records in the group, data stewards may be able to make more informed decisions about what to do with an exception record. For example, in a matching situation a data steward could see all candidates to determine if the exception is a duplicate of the others.
- Note:** If the input data does not contain a field named "CollectionNumber" this option will be disabled.
- Group by** If you selected **Return all records in exception's group**, choose the field by which to group the records.
- Note:** The "CollectionNumber" input field will not appear in this list because it is not a valid selection for the Group by feature.
- Revalidation service** Select the service you want to run when you revalidate records from this dataflow. The service runs when a user saves edited records in the Portal

Exception Editor. Status is changed to Failed for records that fail revalidation. Successfully revalidated records are reprocessed or approved depending on the selection for **Action after revalidation**.

In an approval flow, successfully revalidated records are passed to the next approval level. For the last approval level in an approval flow, revalidated records are either released for reprocessing or retained in the repository as Resolved, depending on the selection for **Action after revalidation**.

Action after revalidation

Specifies whether to reprocess records or approve records that have been successfully revalidated.

- **Reprocess records**—Choose this option to reprocess records that are successfully revalidated. The revalidated records are removed from the repository for reprocessing.
- **Approve records**—Choose this option to approve records that are successfully revalidated. The approved records are retained in the repository and their status changed to Resolved.

Match exception records using match field

Uses match fields to match input records against exception records in the repository. Enable this option if your input contains records that previously generated exceptions but are now corrected in the input.

The input records will be evaluated against the conditions and then matched against the existing exception records in the repository. If an input record passes the conditions and matches an exception record, that exception record will be removed from the repository. If an input record does not pass the conditions and matches an exception record, that exception record will be updated and retained in the repository. Additionally, if duplicates exist in the repository, only one matched exception per dataflow will be updated; all others for that dataflow will be deleted.

Optimized for single records or small batches

This option is activated when you check **Match exception records using match field**. When this option is not checked (default), the server will load into memory all existing exception records for the current dataflow and stage before processing the incoming exception records. This is recommended when the repository has a low number of existing exception records and high number of new exception records or updates. This scenario typically involves a longer initial load time and an increased memory requirement; it is faster when processing larger batches, such as daily, weekly, or monthly updates.

Checking this option is recommended when the repository has a high number of existing exception records and a relatively low number of new exception records or updates, as the server queries the repository for existing exception records as each input record is read in. This scenario typically involves a shorter initial load time and a lower memory requirement; it is faster when processing a few records in real time.

Match fields

Provides a list of all input fields used to build a key to match an exception record in the repository. You must define at least one match field if you checked the **Match exception records using match field** check box.

Add/Modify Condition dialog box

Conditions consist of one or more logical statements that evaluate the value in an input field. Use the **Add Condition** or **Modify Condition** dialog box to define the criteria to determine if a record is an exception.

The **Add Condition** dialog box is opened when you click **Add** on the **Conditions** tab of the **Exception Monitor Options** dialog box. It provides options to configure a new condition. The **Modify Condition** dialog box is opened when you click **Modify**. It allows you to edit values of the same options for an existing condition. After you have create conditions, they appear on the **Conditions** tab of the **Exception Monitor Options** dialog box.

Predefined Conditions Select a predefined condition or retain **<custom condition>** in the drop-down menu to create a new condition. After you have created predefined or custom conditions, they will appear on the **Conditions** tab of the **Exception Monitor Options** dialog box. The icon next to the name of the condition identifies it as either a predefined condition or a custom condition. A dual-document icon  designates a predefined condition, and a single document icon  designates a custom condition. Click the **Save** button to save a condition. After you save a condition, the **Predefined conditions** field changes to show the name of the condition rather than **<custom condition>**.

Name A name for the condition. The name can be anything you like. Since the condition name is displayed in the Data Stewardship Portal, you should use a descriptive name. For example, `MatchScore<80` or `FailedDPV`. If you try to give a new condition a name that is identical to an existing condition but with other characters appended to the end (for example, `FailedDPV` and `FailedDPV2`), you will be asked whether you want to overwrite the existing condition as soon as you type the last character that matches its name (using our example, "V"). Respond **Yes** to the prompt, finish naming the condition, and when you press **OK** or **Save**, both conditions will be visible on the Exception Monitor Options dialog box. The new condition will not overwrite the existing condition unless the new name is identical.

Assign to The user to whom the exception records meeting this condition are assigned. This user is referred to as the *data steward*. If this setting is empty, exception records from a job are automatically assigned to the user who ran the job. For an approval flow, you can configure this setting to specify a data steward other than the user who runs the job.

Condition categories

Data domain (Optional) Specifies the kind of data being evaluated by the condition. This is used solely for reporting purposes to show which types of exceptions occur in your data. For example, if the condition evaluates the success or failure of address validation, the data domain could be "Address"; if the condition evaluates the success or failure of a geocoding operation, the data domain could be "Spatial", and so forth. You can specify your own data domain or select one of the predefined domains:

- Account—The condition checks a business or organization name associated with a sales account.
- Address—The condition checks address data, such as a complete mailing address or a postal code.

- **Asset**—The condition checks data about the property of a company, such as physical property, real estate, human resources, or other assets.
- **Date**—The condition checks date data.
- **Email**—The condition checks email data.
- **Financial**—The condition checks data related to currency, securities, and so forth.
- **Name**—The condition checks personal name data, such as a first name or last name.
- **Phone**—The condition checks phone number data.
- **Product**—The condition checks data about materials, parts, merchandise, and so forth.
- **Spatial**—The condition checks point, polygon, or line data which represents a defined geographic feature, such as flood plains, coastal lines, houses, sales territories, and so forth.
- **SSN**—The condition checks U.S. Social Security Number data.
- **Uncategorized**—Choose this option if you do not want to categorize this condition.

Data quality metric (Optional) Specifies the metric that this condition measures. This is used solely for reporting purposes to show which types of exceptions occur in your data. For example, if the condition is designed to evaluate the record's completeness (meaning, for example, that all addresses contain postal codes) then you could specify "Completeness" as the data quality metric. You can specify your own metric or select one of the predefined metrics:

- **Accuracy**—The condition measures whether the data could be verified against a trusted source. For example, if an address could not be verified using data from the postal authority, it could be considered to be an exception because it is not accurate.
- **Completeness**—The condition measures whether data is missing essential attributes. For example, an address that is missing the postal code, or an account that is missing a contact name.
- **Consistency**—The condition measures whether the data is consistent between multiple systems. For example if your customer data system uses gender codes of M and F, but the data you are processing has gender codes of 0 and 1, the data could be considered to have consistency problems.
- **Interpretability**—The condition measures whether data is correctly parsed into a data structure that can be interpreted by another system. For example, social security numbers should contain only numeric data. If the data contains letters, such as xxx-xx-xxxx, the data could be considered to have interpretability problems.
- **Recency**—The condition measures whether the data is up to date. For example, if an individual moves but the address you have in your system contains the person's old address, the data could be considered to have a recency problem.
- **Uncategorized**—Choose this option if you do not want to categorize this condition.
- **Uniqueness**—The condition measures whether there is duplicate data. If the dataflow could not consolidate duplicate data, the records could be considered to be an exception.

Approval flow

Type Specifies the name of an approval flow type. Approval flows define a succession of review levels through which exception records must be accepted by reviewers after they are edited

by the data steward. When a condition is met, a record is associated with the approval flow specified here. The approval flow type is defined on the Management Console **Resources > Data Stewardship Settings** page. If the option is left as <undefined>, the condition will not be associated with an approval flow type, and records will be resolved by the data steward without any subsequent review.

Expression and Notification tabs

Expressions List expressions defined for a condition. Expressions are logical statements that check the value of a field to determine if the record might be considered an exception. Click the **Add** button to add a new expression. You must add at least one expression to a condition. Expressions are evaluated in the order that they appear here. You can click **Move Up** or **Move Down** to change the order.

Notification Complete options on this tab to send a message to email addresses when this condition is met a specific number of times. A notification email includes a link to the failed records in the Data Stewardship Portal **Editor**, where you can manually enter the correct data. If you do not wish to set up notifications, do not configure options on this tab. To stop sending notifications to a particular email address, remove that address from the list of recipients in the **Send notification to** box.

- **Send notification to**—Specifies email addresses for notifications. You can separate multiple email addresses with semicolons (;) or commas (,).
- **Subject**—Specifies the subject line on email notifications.
- **On**—Specifies to send a notification on the **First occurrence** that an expression evaluates to true or **After** a specified number of occurrences. The maximum allowed value is 1,000,000 occurrences.
- **Send reminder after**—Check this check box and specify the number of days on which to send a repeat email.
- **Remind daily**—Check this check box to continue sending reminders every day until the exception is resolved.
- **Reminder recipients**—Specifies recipients for reminders. The email addresses listed here do not have to match those specified for the original notification.

Add/Modify Expression dialog box

Use the **Add Expression** or **Modify Expression** dialog box to define expressions, which are logical statements that check the value of a field to determine if the record might be considered an exception. You must add at least one expression to a condition.

This dialog box opens when you click **Add** or **Modify** in the **Add/Modify Condition** dialog box.

Expression type

- | | |
|---|--|
| Expression created with Expression Builder | Select this option to create a basic expression. |
| Custom expression | Select this option to write an expression using Groovy scripting. If you need to use more complex logic, such as nested evaluations, |

use a custom expression. For more information, see [Using Custom Expressions in Exception Monitor](#) on page 318.

Expression

Logical operator If other expressions are already defined for this condition, you can select an operator to combine it with the preceding expressions.

- **And**—This expression must be true in addition to the preceding expression being true in order for the condition to be true.
- **Or**—If this expression is true the condition is true even if the preceding expression is not true.

If you chose **Expression created with expression builder** for the **Expression type**, the following options are available.

Field name	Select the field that you want this expression to evaluate. The list of available fields is populated based on the stages upstream from the Exception Monitor stage.
Operator	Select the operator you want to use in the evaluation.
Value	Specify the value you want the expression to check for using the operator chosen in the Operator field. This is a regular expression if you selected matches regular expression as the Operator .

How to

Add the Exception Monitor stage to a workflow

1. In the **Stages** palette, expand **Primary Stages > deployed Stages > Data Stewardship**.
2. Drag the **Exception Monitor** stage to the canvas.
3. Optional: Click the Exception Monitor label and type a meaningful name for the stage in the workflow.
4. Drag a connection from an upstream stage to the input port on the Exception Monitor stage. This is the path through which records will enter the Exception Monitor.
5. Drag a connection from the success port to the first downstream stage that will process records that do not trigger any exceptions.
6. Drag a connection from the exception port to the first downstream stage that will process records that trigger an exception.
7. Double-click the Exception Monitor stage to configure its options.
8. On the **Conditions** tab, click the **Add** button to add conditions in the order that you want them evaluated.
 - a) In the **Add Condition** dialog box, complete **Name**, **Assign to**, **Data domain**, and **Data quality metric**. Under **Approval flow**.
For more information, see [Add/Modify Condition dialog box](#) on page 314.

- b) Optional: In the **Type** box under **Approval flow**, you can select an approval type to trigger an approval flow with a condition.
To show on the **Conditions** tab, approval types must be defined on the Management Console **Data Stewardship Settings** page.
- c) On the **Expressions** tab, click **Add** to add expressions to the condition.
You must add at least one expression to a condition.
For more information, see [Expression and Notification tabs](#) on page 316.
- d) After you add expressions, click the **Move Up** or **Move Down** buttons to change the order in which expressions are evaluated.
- e) Optional: On the **Notification** tab, you can add recipients for reminder emails.
For more information, see [Expression and Notification tabs](#) on page 316.
- f) Click **OK**.
- g) To save this condition for reuse as a predefined condition, click **Save**.
For an existing condition, you will be prompted whether to overwrite the condition.

Note: If you overwrite a predefined condition, any changes will affect all dataflows that already use the condition.

- 9. Use the **Move Up** and **Move Down** buttons to change the order in which conditions are evaluated.
The order of the conditions is important only if you have enabled the option **Stop evaluating when a condition is met**. For more information, see [Conditions tab](#) on page 311
- 10. Configure options on the **Configuration** tab.
For more information, see [Configuration tab](#) on page 312.
- 11. Click the **OK** button.

Using Custom Expressions in Exception Monitor

You can write your own custom expressions to control how Exception Monitor routes records using the Groovy scripting language to create an expression.

Using Groovy Scripting

For information on Groovy, see groovy-lang.org.

Groovy expressions used in the Exception Monitor stage must evaluate to a Boolean value (true or false) that indicates whether the record is considered an exception and should be routed for manual review. Exception records are routed to the exception port.

For example, if you need to review records with a validation confidence level of <85, your script would look like:

```
data['Confidence']<85
```

The monitor would evaluate the value of the Confidence field against your criteria to determine which output port to send it to.

Checking a Field for a Single Value

This example evaluates to true if the Status field has 'F' in it. This would have to be an exact match, so 'f' would not evaluate to true.

```
return data['Status'] == 'F';
```

Checking a Field for Multiple Values

This example evaluates to true if the Status field has 'F' or 'f' in it.

```
boolean returnValue = false;
if (data['Status'] == 'F' || data['Status'] == 'f')
{
    returnValue = true;
}
return returnValue;
```

Evaluating Field Length

This example evaluates to true if the PostalCode field has more than 5 characters.

```
return data['PostalCode'].length() > 5;
```

Checking for a Character Within a Field Value

This example evaluates to true if the PostalCode field has a dash in it.

```
boolean returnValue = false;
if (data['PostalCode'].indexOf('-') != -1)
{
    returnValue = true;
}
return returnValue;
```

Common Mistakes

The following illustrate common mistakes when using scripting.

The following is incorrect because PostalCode (the column name) must be in single or double quotes

```
return data[PostalCode];
```

The following is incorrect because no column is specified

```
return data[];
```

The following is incorrect because `row.set()` does not return a Boolean value. It will always evaluate to false as well as change the PostalCode field to 88989.

```
return row.set('PostalCode', '88989');
```

Use a single equals sign to set the value of a field, and a double equals sign to check the value of a field.

Read Exceptions

Read Exceptions is a stage that reads records from the exception repository as input to a dataflow. (For more information on the exception repository, see [Data Stewardship Portal](#).)

Note: Once a record is read into a dataflow by Read Exceptions, it is deleted from the repository.

Input

The Read Exceptions stage reads data in from an exception repository. It does not take input from another stage in a dataflow.

Note: Only records marked as Resolved in the Data Stewardship Portal are read into the dataflow. This option can be changed by clearing the **Process resolved records** check box on the **Runtime** tab.

Output

The Read Exceptions stage has a primary Output port and an optional History port.

Output port

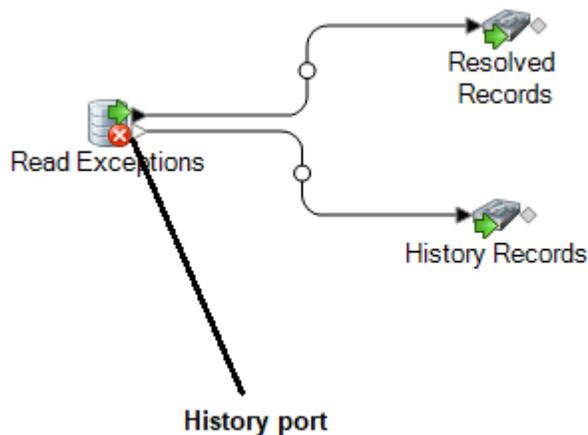
The primary Output port returns records from the Data Stewardship repository that have resolved status and that match the selection criteria specified in the Read Exception options. In addition to the record fields, Read Exceptions returns the following fields that describe the last modifications made to the record in the Data Stewardship Portal.

Table 22: Read Exceptions Output

Field Name	Description
Exception.Comment	Any comments entered by the person who resolved the exception. For example, comments might describe the modifications that the Data Stewardship made to the record.
Exception.LastModifiedBy	The last user to modify the record in the Data Stewardship Portal
Exception.LastModifiedMilliseconds	The time that the record was last modified in the Data Stewardship Portal. The time is expressed in milliseconds since January 1, 1970 0:00 GMT. This is the standard way of calculating time in the Java programming language. You can use this value to perform date comparisons, or to create a transform to convert this value to whatever date format you want.
Exception.LastModifiedString	The time that the record was last modified in the Data Stewardship Portal. This field provides a more understandable representation of the date than the Exception.LastModifiedMilliseconds field. The time is expressed in this format: Thu Feb 17 13:34:32 CST 2011

History port

The optional history port streams all versions of records. History records are typically used to assemble change logs.



Options

Settings in the Read Exceptions stage are configured on the **General**, **Sort**, and **Runtime** tabs.

General Tab

The options on the **General** tab specify which exception records you want to read into the dataflow.

The **Filter** options allow you to select a subset of records from the exception repository using these criteria.

User	The user who ran the dataflow that generated the exceptions you want to read into the dataflow.
Dataflow name	The name of the dataflow that generated the exceptions you want to read into the dataflow.
Stage label	The Exception Monitor stage's label as shown in the dataflow in Enterprise Designer. This criteria is useful if the dataflow that generated the exceptions contains multiple Exception Monitor stages and you only want to read in the exceptions from one of those Exception Monitor stages.
From date	The date and time of the oldest records that you want to read into the dataflow. The date of an exception record is the date it was last modified.
To date	The date and time of the newest records that you want to read into the dataflow. The date of an exception record is the date it was last modified.

The **Fields** listing shows the fields that will be read into the dataflow. By default all fields are included, but you can exclude fields by clearing the check box in the **Include** column.

The **Preview** listing shows the records that meet the criteria you specified under **Filter**.

Note: The preview displays only records that have been marked "Approved" in the Data Stewardship Portal and meet the filter criteria.

Sort Tab

Options on the **Sort** tab specify how to sort input records based on field values.

Add	Adds a field to sort on.
Field Name column	Shows the name of the field to sort on. You can select a field by clicking the drop-down button.
Order column	Specifies whether to sort in ascending or descending order.
Up and Down	Changes the order of the sort. Records are sorted first by the field at the top of the list, then by the second, and so forth.
Remove	Removes a sort field.

Runtime Tab

Options on the **Runtime** tab specify runtime behavior of the Read Exceptions stage.

Starting record	Specify the position in the repository of the first record you want to read into the dataflow. For example, if you want to skip the first 99 records in the repository, you would specify 100. The 100th record would be the first one read into the repository if it matches the criteria specified on the General tab. A record's position is determined by the order of the records in the Data Stewardship Portal.
All records	Select this option if you want to read in all records that match the search criteria specified on the General tab.
Max records	Select this option if you want to limit the number of records read in to the dataflow. For example, if you only want to read in the first 1,000 records that match the selection criteria, select this option and specify 1000.
Process resolved records	Select this check box to process resolved records. Clear this check box to process and purge unresolved records. By default, this check box is selected.

Write Exceptions

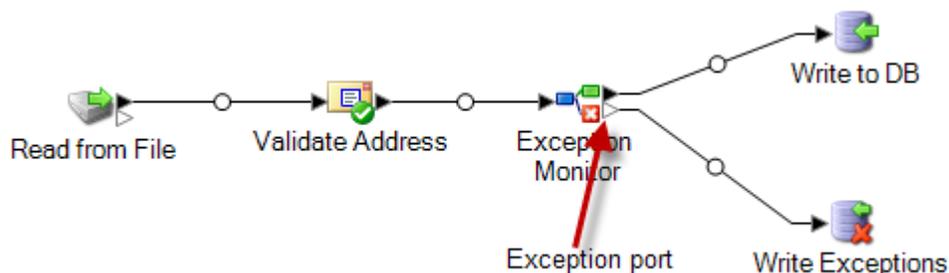
Write Exceptions is a stage that takes records that the Exception Monitor stage has identified as exceptions and writes them to the exception repository. Once in the exception repository, the records can be reviewed and edited using the Data Stewardship Portal.

Note: Exception records are not written to the Data Stewardship Portal when jobs or services are run in inspection mode in Enterprise Designer or preview mode in Management Console.

Input

The Write Exceptions stage takes records from the exception port on the Exception Monitor stage and then writes them to the exception repository. The Write Exceptions stage should be placed

downstream of the Exception Monitor stage's exception port. The exception port is the bottom output port on the Exception Monitor stage:



Output

Write Exceptions does not return any output in the dataflow. It writes exception records to the exception repository.

Options

The Write Exceptions stage enables you to select which fields' data should be returned to the exceptions repository. The fields that appear depend upon the stages that occur upstream in the dataflow. If, for instance, you have a Validate Address stage in the dataflow, you would see such fields as AddressLine1, AddressLine2, City, PostalCode, and so on in the Write Exceptions stage. By default, all of those fields are selected; uncheck the boxes for any fields you do not want returned to the exceptions repository. The order of the fields is determined by how they are ordered when they come into the Write Exceptions stage. You can reorder the fields by selecting a row and using the arrows on the right side of the screen to move the row up or down. The order you select here will persist for all users in the Data Stewardship Portal, but each user can reorder the fields within the Portal to their own liking.

Select **Allow user to create best of breed records in Portal** to perform manual record consolidation when data matches cannot be made using standard rules. This feature copies a selected record within a group and uses it instead of the duplicates for processing. This option is available only for grouped exception records that are generated from a matching job, which is identified by the presence of the CandidateGroup field or the CollectionNumber field. When you use this option, a read-only field called "CollectionRecordType" will be added to the exception record. You can see this field at the bottom of the list; note that all options for that field are disabled.

Note: If your dataflow is also being configured for **revalidation**, you will need to manually add and expose the CollectionRecordType field in the Exception Monitor stage/subflow and the service itself.

After adding a **Best of Breed** record in the Resolve Duplicates view of the Data Stewardship Portal **Editor**, this field will be set to "BestOfBreed." If you choose to create best of breed records, you will

be unable to use the **Approve All** option for those records in the Data Stewardship Portal. Read more about best of breed records in the Data Stewardship Portal [here](#).

You may have input fields that you want in the repository but do not want to be viewable within the Data Stewardship Portal. This could be due to the field containing sensitive data or simply because you want to streamline what appears in the Portal. Check the **Allow viewing** box to designate which of the selected fields should be viewable once they are passed to the exceptions repository. By default, all fields are viewable. Uncheck the box for any field you do not want visible in the Portal.

Additionally, you can designate which of the selected fields should be editable in the Portal once they are passed to the exceptions repository. By default, the **Allow editing** column is checked for all fields coming in to the Write Exceptions stage. Uncheck the box for any field you wish to be returned to the exceptions repository in a read-only state.

Finally, you can use the **Lookup** function to assign a lookup to a field containing problematic data. You can select from the list of lookups that have been defined in the Data Stewardship Settings tool or you can manually enter the name of the lookup. For more information on lookups, see [Lookups](#) on page 236.

Note: Lookups can be assigned only to fields whose type is string.

Data Normalization

Spectrum Data Normalization

Spectrum Data Normalization examines terms in a record and determines if the term is in the preferred form.

Components

Spectrum Data Normalization consists of:

- **Advanced Transformer**—This stage scans and splits strings of data into multiple fields, placing the extracted and non-extracted data into a new or existing field.
- **Open Parser**—This stage parses your input data from many cultures of the world using a simple, but powerful parsing grammar. Using this grammar, you can define a sequence of expressions that represent domain patterns for parsing your input data. Open Parser also collects statistical data and scores the parsing matches to help you determine the effectiveness of your parsing grammars.

- **Table Lookup**—This stage evaluates a term and compares it to a previously validated form of that term. If the term is not in the proper form, then the standard version replaces the term. Table Lookup can change full words to abbreviations, change abbreviations to full words, change nicknames to full names, and can correct misspellings.
- **Transliterator**—Transliterator converts a string between Latin and other scripts.

Advanced Transformer

The Advanced Transformer job scans and splits strings of data into multiple fields using tables or regular expressions. It extracts a specific term or a specified number of words to the right or left of a term. Extracted and non-extracted data can be placed into an existing field or a new field.

For example, want to extract the suite information from this address field and place it in a separate field.

2300 BIRCH RD STE 100

To accomplish this, you could create an Advanced Transformer that extracts the term STE and all words to the right of the term STE, leaving the field as:

2300 BIRCH RD

Input

Advanced Transformer uses any defined input field in the data flow.

Options

Advanced Transformer options can be configured at the stage level, through any of the Spectrum Technology Platform clients, or at runtime, using dataflow options.

Configuring Options

To specify the options for Advanced Transformer you create a rule. You can create multiple rules then specify the order in which you want to apply the rules. To create a rule:

1. Double-click the instance of Advanced Transformer on the canvas. The Advanced Transformer Options dialog displays.
2. Select the number of runtime instances and click **OK**. Use the Runtime Instances option to configure a dataflow to run multiple, parallel instances of a stage to potentially increase performance.
3. Click the **Add** button. The Advanced Transformer Rule Options dialog displays.

Note: If you add multiple transformer rules, you can use the **Move Up** and **Move Down** buttons to change the order in which the rules are applied.

4. Select the type of transform action you wish to perform and click **OK**. The options are listed in the table below.

Table 23: Advanced Transformer Options

Option	Description
Source	Specifies the source input field to evaluate for scan and split.
Extract using	<p>Select Table Data or Regular Expressions.</p> <p>Select Table Data if you want to scan and split using the XML tables located at <Drive>:\Program Files\Precisely\Spectrum\server\modules\advancedtransformer\data. See Table Data Options below for more information about each option.</p> <p>Select Regular Expressions if you want to scan and split using regular expressions. Regular expressions provide many additional options for splitting data. You can use the prepackaged regular expressions by selecting one from the list or you can construct your own using RegEx syntax.</p> <p>For example, you could split data when the first numeric value is found, as in "John Smith 123 Main St." where "John Smith" would go in one field and "123 Main St." would go in another. See Regular Expression options below for more information about each option.</p>
Table Data Options	
Non-extracted Data	<p>Specifies the output field that you want to contain the transformed data. If you want to replace the original value specify the same field in the Destination field as you did in the Source drop-down box.</p> <p>You may also type in a new field name in the Destination field. If you type in a new field name, that field name will be available in stages in your dataflow that are downstream of Advanced Transformer.</p>
Extracted Data	<p>Specifies the output field where you want to put the extracted data.</p> <p>You may type in a new field name in the Extracted Data field. If you type in a new field name, that field name will be available in stages in your dataflow that are downstream of Advanced Transformer.</p>

Option	Description
Tokenization Characters	<p>Specifies any special characters that you want to tokenize. Tokenization is the process of separating terms. For example, if you have a field with the data "Smith, John" you would want to tokenize the comma. This would result in terms:</p> <ul style="list-style-type: none"> • Smith • , • John <p>Now that the terms are separated, the data can be split by scanning and extracting on the comma so that "Smith" and "John" are cleanly identified as the data to standardize.</p>
Table	<p>Specifies the table that contains the terms on which to base the splitting of the field. For a list of tables, see Advanced Transformer Tables on page 150. For information about creating or modifying tables, see Introduction to Lookup Tables on page 150.</p>
Lookup multiple word terms	<p>Select this check box to enable multiple word searches within a given string. For example:</p> <p>Input String = "Cedar Rapids 52401" Business Rule = Identify "Cedar Rapids" in string based on a table that contains the entry; Cedar Rapids = US Output = Identifies presence of "Cedar Rapids" and places the terms into a new field, for example City.</p> <p>For multiple word searches, the search stops at the first occurrence of a match.</p> <p>Note: Selecting this option may adversely affect performance.</p>

Option	Description
Extract	<p>Specifies the type of extraction to perform. Select from one of these:</p> <p>Extract term Extracts the term identified by the selected table.</p> <p>Extract N words to the right of the term Extracts words to the right of the term. You specify the number of words to extract. For example, if you want to extract the two words to the right of the identified term, specify 2.</p> <p>Extract N words to the left of the term Extracts words to the left of the term. You specify the number of words to extract. For example, if you want to extract the two words to the left of the identified term, specify 2.</p> <p>If you choose to extract words to the right or left of the term, you can specify if you want to include the term itself in the destination data or the extracted data. For example, if you have this field:</p> <p>2300 BIRCH RD STE 100</p> <p>and you want to extract "STE 100" and place it in the field specified in extracted data, you would choose to include the term in the extracted data field, thus including the abbreviation "STE" and the word "100".</p> <p>If you select neither Destination nor Extracted data, the term will not be included and is discarded.</p>
Regular Expressions Options	
Regular Expressions	<p>Select a prepackaged regular expressions from the list or construct your own in the text box. Advanced Transformer supports standard RegEx syntax.</p> <p>The Java 2 Platform contains a package called java.util.regex, enabling the use of regular expressions. For more information, go to: java.sun.com/docs/books/tutorial/essential/regex/index.html.</p>
Ellipsis Button	Click this button to add or remove a new regular expression.
Populate Group	After you have selected a predefined or typed a new Regex expression, click Populate Group to extract any Regex groups and place the complete expression, as well as any Regex groups found, into the Groups list.

Option	Description
Groups	<p>This column shows the regular expressions for the selected Regular Expressions group.</p> <p>For example, if you select the Date Regex expression, the following expression displays in the text box: (1[012]{1,2}0?[1-9])[-./]([12][0-9])3[01]{1,2}0?[1-9])[-./]([0-9]{4}). This Regex expression has three parts to it and the whole expression and each of the parts can be sent to a different output field. The entire expression is looked for in the source field and if a match is found in the source field, then the associated parts are moved to the assigned output field. If the source field is "On 12/14/2006" and you apply the Date expression to it, and assign the entire date (such as, "12/14/2006) to be placed in the DATE field, the "12" to be placed in MONTH field, the "14" to be placed in the DAY field and "2006" to be placed in YEAR field. It will look for the date and if it finds it will move the appropriate information to the appropriate output field.</p> <p>Source Field: "On 12/14/2006" DATE: "12/14/2006" MONTH: "12" DAY: "14" YEAR: "2006"</p>
Output Field	Pull-down menu to select an output field.

Configuring Options at Runtime

Advanced Transformer rules can be configured and passed at runtime if they are exposed as dataflow options. This enables you to override the existing configuration with JSON-formatted strings. You can also set stage options when calling the job through a process flow or through the job executor command-line tool.

You can find schemas for AdvancedTransformerRules in the following folder:

```
<Spectrum Location>\server\modules\jsonSchemas\advancedTransformer
```

To define Advanced Transformer rules at runtime:

1. In Enterprise Designer, open a dataflow that uses the Advanced Transformer stage.
2. Save and expose that dataflow.
3. Go to **Edit > Dataflow Options**.
4. In the **Map dataflow options to stages** table, expand Advanced Transformer. Check the box for AdvancedTransformerRules.
5. Optional: Change the name of the options in the **Option label** field.
6. Click **OK** twice.

Output

Advanced Transformer does not create any new output fields. Only the fields you define are written to the output.

Open Parser

Open Parser parses your input data from many cultures of the world using a simple but powerful parsing grammar. Using this grammar, you can define a sequence of expressions that represent domain patterns for parsing your input data. Open Parser also collects statistical data and scores the parsing matches to help you determine the effectiveness of your parsing grammars.

Use Open Parser to:

- Parse input data using domain-specific and culture-specific parsing grammars that you define in Domain Editor.
- Parse input data using domain-independent parsing grammars that you define in Open Parser using the same simple but powerful parsing grammar available in Domain Editor.
- Parse input data using domain-independent parsing grammars at runtime that you define in Dataflow Options.
- Preview parsing grammars to test how sample input data parses before running the job using the target input data file.
- Trace parsing grammar results to view how tokens matched or did not match the expressions you defined and to better understand the matching process.

Input

Open Parser accepts the input fields that you define in your parser grammar. For more information, see [Header Section Commands](#) on page 30.

If you are performing culture-specific parsing, you can optionally include a CultureCode field in the input data to use a specific culture's parsing grammar for a record. If you omit the CultureCode field, or if it is empty, then each culture listed in the Open Parser stage is applied, in the order specified. The result from the culture with the highest parser score, or the first culture to have a score of 100, is returned. For more information about the CultureCode field, see [Assigning a Parsing Culture to a Record](#) on page 12.

Options

The following tables list the options for the Open Parser stage.

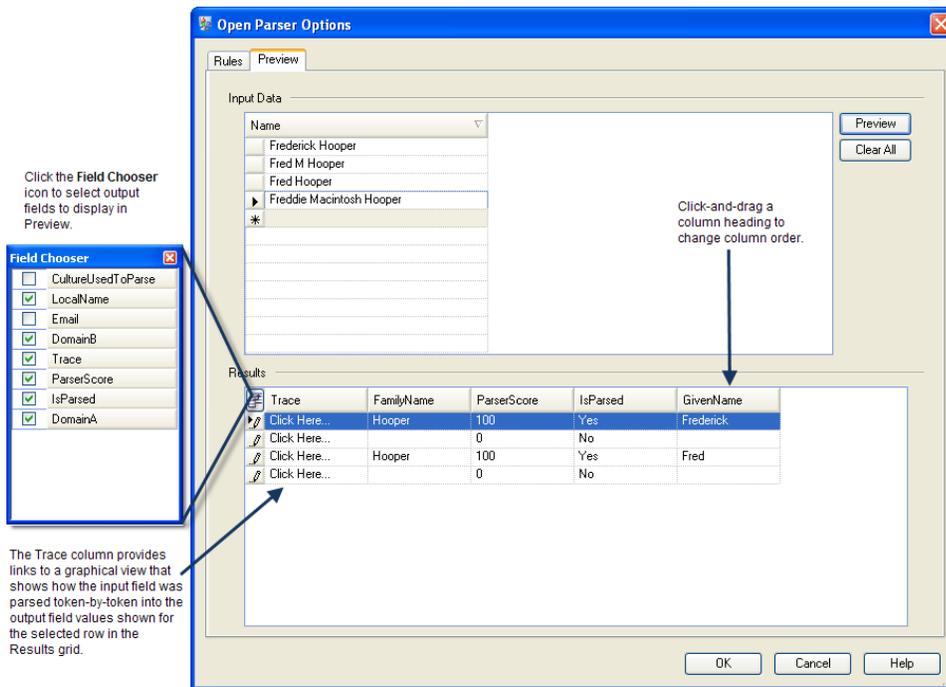
Rules Tab

Option	Description
Use culture-specific domain grammar	<p>Specifies to use a language and domain specific parsing grammar which has already been defined in the Open Parser Domain Editor tool in Enterprise Designer. For more information about defining domains, see Defining a Culture-Specific Parsing Grammar on page 10.</p> <p>If you choose this option you will also see these options:</p> <p>Domain Specifies the parsing grammar to use.</p> <p>Cultures Specifies the language or culture of the data you want to parse. Click the Add button to add a culture. You can change the order in which Open Parser attempts to parse the data with each culture by using the Move Up and Move Down buttons. For more information about cultures, see Defining a Culture-Specific Parsing Grammar on page 10.</p> <p>Return multiple parsed records Enable this option to have Open Parser return records for each culture that successfully parses the input. If you do not check this box, Open Parser will return the results for the first record that achieves a parser score of 100, regardless of culture. If all cultures run without hitting a record that has parser score of 100, Open Parser will return the record with the score closest to 100. If multiple cultures return records with the same high score under 100, the order set in Step 4 will determine which culture's record is returned.</p>
Define domain-independent grammar	<p>Choose this option if you want to define a parsing grammar that should be applied without consideration of the language or domain of the input data. If you choose this option, the grammar editor will appear and you can define the parsing grammar directly in the Open Parser stage rather than using the Open Parser Domain Editor tool in Enterprise Designer.</p> <p>Note: You can also define domain-independent grammar at runtime. For more information, see Defining Domain-Independent Parsing Grammars at Runtime.</p>

Preview Tab

Creating a working parsing grammar is an iterative process. Preview is useful in testing out variations on your input to make sure that the parsing grammar produces the expected results.

Type test values in the input field and then click **Preview**.



The parsed output fields display in the **Results** grid. For information about the output fields, see [Output](#) on page 333. For information about trace, see [Tracing Final Parsing Results](#) on page 36. If your results are not what you expected, click the **Rules** tab and continue editing the parsing grammar and testing input data until it produces the expected results.

Output

Table 24: Open Parser Output

Field Name	Description / Valid Values
<Input Field>	The original input field defined in the parsing grammar.
<Output Fields...>	The output fields defined in the parsing grammar.
CultureCode	The culture codes contained in the input data. For a complete list of supported culture codes, see Assigning a Parsing Culture to a Record on page 12.
CultureUsedtoParseSelect a match	The culture code value used to parse each output record. This value is based on results in the Match Results List and matches to a culture-specific parsing grammar. then click Remove .

Field Name	Description / Valid Values
IsParsed	Indicates whether an output record was parsed. The possible values are Yes or No.
ParserScore Select a match results in the Match Results List and then click Remove .	Indicates the total average score. The value of ParserScore will be between 0 and 100, as defined in the parsing grammar. 0 is returned when no matches are returned. For more information, see Scoring .
Trace	Click this control to see a graphical view of how each token in the parsing grammar was parsed to an output field for the selected row in the Results grid.

Table Lookup

The Table Lookup stage standardizes terms against a previously validated form of that term and applies the standard version. This evaluation is done by searching a table for the term to standardize.

For example:

	First Name	Last Name
Source Input:	Bill	Smith
Standardized Output:	William	Smith

There are three types of action you can perform: standardize, identify, and categorize.

If the term is found when performing the standardize action, Table Lookup replaces either the entire field or individual terms within the field with the standardized term, even if the field contains multiple words. Table Lookup can include changing full words to abbreviations, changing abbreviations to full words, changing nicknames to full names or misspellings to corrected spellings.

If the term is found when performing the identify action, Table Lookup flags the record as containing a term that can be standardized, but performs no action.

If the term is found when performing the categorize action, Table Lookup uses the source value as a key and copies the corresponding value from the table entry into the selected field. If none of the source terms match, **Categorize** uses the default value specified.

Input

Table 25: Table Lookup Input Fields

Field Name	Description / Valid Values
Source	Specifies the source input field to evaluate for scan and split.
StandardizationTable	One of the tables listed in Table Lookup Tables on page 153.

Options

Table Lookup options can be configured at the stage level, through any of the Spectrum Technology Platform clients, or at runtime, using dataflow options.

Configuring Options

To specify the options for Table Lookup you create a rule. You can create multiple rules then specify the order in which you want to apply the rules. To create a rule, open the Table Lookup stage and click **Add** then complete the following fields.

Note: If you add multiple Table Lookup rules, you can use the **Move Up** and **Move Down** buttons to change the order in which the rules are applied.

Option	Description
Action	<p>Specifies the type of action to take on the source field. One of the following:</p> <p>Standardize Changes the data in a field to match the standardized term found in the lookup table. If the field contains multiple terms, only the terms that are found in the lookup table are replaced with the standardized term. The other data in the field is not changed.</p> <p>Identify Flags the record as containing a term that can be standardized, but performs no action on the data in the field. The output field StandardizedTermIdentified is added to the record with a value of Yes if the field can be standardized and No if it cannot.</p> <p>Categorize Uses the Source value as a key and copies the corresponding value from the table into the field selected in the Destination list. This creates a new field in your data that can be used to categorize records.</p>
On	<p>Specifies whether to use the entire field as the lookup term or to search the lookup table for each term in the field. One of the following:</p> <p>Complete field Treats the entire field as one term, resulting in the following:</p> <ul style="list-style-type: none"> • If you selected the action Standardize, Table Lookup treats the entire field as one string and attempts to standardize the field using the string as a whole. For example, "International Business Machines" would be changed to "IBM". • If you selected the action Identify, Table Lookup treats the entire field as one string and flags the record if the string as a whole can be standardized. • If you selected the action Categorize, Table Lookup treats the entire field as one string and flags the record if the string as a whole can be categorized. <p>Individual terms within field Treats each word in the field as its own term, resulting in the following:</p> <ul style="list-style-type: none"> • If you selected the action Standardize, Table Lookup parses the field and attempts to standardize the individual terms within the field. For example, "Bill Mike Smith" would be changed to "William Michael Smith." • If you selected the action Identify, Table Lookup parses the field and flags the record if any single term within the field can be standardized. • If you selected the action Categorize, Unlike Standardize, Categorize does not copy the source term if there isn't a table match. If none of the source terms match, Categorize uses the default value specified. Unlike Standardize, Categorize only returns that table value and nothing from Source. If none of the source terms match, Categorize uses the default value specified.

Option	Description				
Source	Specifies the field you want to containing the term you want to look up.				
Destination	<p>Specifies the field to which the terms returned by the table lookup should be written.</p> <p>If you want to replace the value, specify the same field in the Destination field as you did in the Source field. You can also create a new field by typing the name of the field you want to create.</p> <p>The Destination field is not available if you select the action Identify.</p>				
Table	<p>Specifies the table you want to use to find terms that match the data in your dataflow.</p> <p>For a list of tables that you can edit, see Table Lookup Tables on page 153. For information about creating or modifying tables, see Introduction to Lookup Tables on page 150.</p>				
Lookup multiple word terms	<p>Enables multiple word searches within a given string. For example:</p> <p>Input String: "Major General John Smith"</p> <p>Business Rule: Identify "Major General" in a string based on a table that contains the entry</p> <p>Output: Replace "Major General" with "Maj. Gen."</p> <p>For multiple word searches, the search stops at the first occurrence of a match.</p> <p>This option is disabled when On is set to Complete field.</p> <p>Note: Selecting this option may adversely affect performance.</p>				
When table entry not found, set Destination's value to	<p>Specifies the value to put in the destination field if a matching term cannot be found in the lookup table. One of the following:</p> <table border="0"> <tr> <td style="padding-right: 20px;">Source's value</td> <td>Put the value from the source field into the destination field.</td> </tr> <tr> <td>Other</td> <td>Put a specific value into the destination field.</td> </tr> </table>	Source's value	Put the value from the source field into the destination field.	Other	Put a specific value into the destination field.
Source's value	Put the value from the source field into the destination field.				
Other	Put a specific value into the destination field.				

Configuring Options at Runtime

Table Lookup options can be configured and passed at runtime if they are exposed as dataflow options. This enables you to override the existing configuration with JSON-formatted strings. You can also set stage options when calling the job through a process flow or through the job executor command-line tool.

You can find a schema for LookupRule in the following folder:

```
<Spectrum Location>\server\modules\jsonSchemas\tableLookup
```

To define Table Lookup rules at runtime:

1. In Enterprise Designer, open a dataflow that uses the Table Lookup stage.
2. Save and expose that dataflow.
3. Go to `Edit > Dataflow Options`.
4. In the **Map dataflow options to stages** table, expand Table Lookup. Check the box for LookupRule.
5. Optional: Change the name of the options in the **Option label** field.
6. Click **OK** twice.

Output

Table 26: Table Lookup Outputs

Field Name	Description / Valid Values
StandardizedTermIdentified	Indicates whether or not the field contains a term that can be standardized. Only output if you select Complete field or Individual terms in field options.
	Yes The record contains a term that can be standardized.
	No The record does not contain a term that can be standardized.

Transliterator

Transliterator converts a string between Latin and other scripts. For example:

Source	Transliteration
キャンパス	kyanpasu
Αλφαβητικός Κατάλογος	

Source	Transliteration
биологическом	biologichyeskom

It is important to note that transliteration is not translation. Rather, transliteration is the conversion of letters from one script to another without translating the underlying words.

Note: Standard transliteration methods often do not follow the pronunciation rules of any particular language in the target script.

The Transliterator stage supports these scripts. In general, the Transliterator stage follows the UNGEGN Working Group on Romanization Systems guidelines. For more information, see www.eki.ee/wgrs.

Arabic	The script used by several Asian and African languages, including Arabic, Persian, and Urdu.
Cyrillic	The script used by Eastern European and Asian languages, including Slavic languages such as Russian. The Transliterator stage generally follows ISO 9 for the base Cyrillic set.
Devanagari	The script used by several Indian languages, including Hindi and Sanskrit. This script is a descendent of the Brahmi script which is one of the oldest writing systems used in Ancient India and present South and Central Asia.
Greek	The script used by the Greek language. This script belongs to the Hellenic branch of the Indo-European language family.
Gujarati	The script used by the state of Gujarat in western India. It is one of the modern scripts of India which was adapted from the Devanagari script.
Gurmukhi	The script used by Indian language Punjabi. This script has a considerable influence from Nagari script which is an earlier form of the Devanagari script.
Hangul	The script used by the Korean language. The Transliterator stage follows the Korean Ministry of Culture and Tourism Transliteration regulations. For more information, see the website of The National Institute of the Korean Language .
Han	The script used by Chinese language. It is a branch of the Tibetan-Burman language family and has been written with scripts based on Thai and Chinese.
Traditional/Simplified Chinese	The Transliterator stage supports both traditional and simplified Chinese. For example, this is Traditional Chinese: 人人生而自由. This is Simplified Chinese: 人人生而自由
Kannada	The script used by several South Indian languages, such as Konkani. This script is a descendent of Brahmi script of ancient India.

Katakana and Hiragana	One of several scripts that can be used to write Japanese. The Transliterator stage uses a slight variant of the Hepburn system. With Hepburn system, both ZI (ジ) and DI (ヂ) are represented by "ji" and both ZU (ズ) and DU (ヅ) are represented by "zu". This is amended slightly for reversibility by using "dji" for DI and "dzu" for DU. The Katakana transliteration is reversible. Hiragana-Katakana transliteration is not completely reversible since there are several Katakana letters that do not have corresponding Hiragana equivalents. Also, the length mark is not used with Hiragana. The Hiragana-Latin transliteration is also not reversible since internally it is a combination of Katakana-Hiragana and Hiragana-Latin.
Half width/Full width	The Transliterator stage can convert between narrow half-width scripts and wider full-width scripts. For example, this is half-width: アルアノリウ. This is full-width: アルアノリウ.
Latin	The script used by most languages of Europe, such as English. It was originally used by the ancient Romans to write the Latin language.
Malayalam	The script used by the Malayalam language, the official language of the Indian state of Kerala. This script was first written with the Vatteluttu alphabet which means 'round writing' and developed from the Brahmi script of ancient India.
Oriya	The script used by the Oriya language, the official language of the Indian state of Odisha. The Oriya script was developed from the Kalinga script, one of the many descendents of the Brahmi script of ancient India.
Tamil	The script used by the Tamil language in several states of India, Sri Lanka, and Malaysia. This script was originally written with a version of the Brahmi script known as Tamil Brahmi.
Telugu	The script used by several languages of South India. This script is a descendent of Brahmi script of ancient India.
Thai	The script used by Thai language. This script is influenced by the Brahmi script of ancient India and the Khmer alphabets.

Transliterator is a part of Data Normalization. For a listing of other stages, see [Spectrum Data Normalization](#) on page 325.

Transliteration Concepts

There are a number of generally desirable qualities for script transliterations. A good transliteration should be:

- Complete
- Predictable
- Pronounceable
- Unambiguous

These qualities are rarely satisfied simultaneously, so the Transliterator stage attempts to balance these requirements.

Complete

Every well-formed sequence of characters in the source script should transliterate to a sequence of characters from the target script.

Predictable

The letters themselves (without any knowledge of the languages written in that script) should be sufficient for the transliteration, based on a relatively small number of rules. This allows the transliteration to be performed mechanically.

Pronounceable

Transliteration is not as useful if the process simply maps the characters without any regard to their pronunciation. Simply mapping "αβγδεζηθ..." to "abcdefgh..." would yield strings that might be complete and unambiguous, but cannot be pronounced.

Standard transliteration methods often do not follow the pronunciation rules of any particular language in the target script. For example, the Japanese Hepburn system uses a "j" that has the English phonetic value (as opposed to French, German, or Spanish), but uses vowels that do not have the standard English sounds. A transliteration method might also require some special knowledge to have the correct pronunciation. For example, in the Japanese kunrei-siki system, "tu" is pronounced as "tsu". This is similar to situations where there are different languages within the same script. For example, knowing that the word Gewalt comes from German allows a knowledgeable reader to pronounce the "w" as a "v".

In some cases, transliteration may be heavily influenced by tradition. For example, the modern Greek letter beta (β) sounds like a "v", but a transform may continue to use a b (as in biology). In that case, the user would need to know that a "b" in the transliterated word corresponded to beta (β) and is to be pronounced as a "v" in modern Greek. Letters may also be transliterated differently according to their context to make the pronunciation more predictable. For example, since the Greek sequence GAMMA GAMMA (γγ) is pronounced as "ng", the first GAMMA can be transcribed as an "n".

Note: In general, in order to produce predictable results when transliterating Latin script to other scripts, English text will not produce phonetic results. This is because the pronunciation of English cannot be predicted easily from the letters in a word. For example, grove, move, and love all end with "ove", but are pronounced very differently.

Unambiguous

It should always be possible to recover the text in the source script from the transliteration in the target script. For example, it should be possible to go from Elláda back to the original Ελλάδα. However, in transliteration multiple characters can produce ambiguities. For example, the Greek

character PSI (ψ) maps to ps, but ps could also result from the sequence PI, SIGMA ($\pi\sigma$) since PI (π) maps to p and SIGMA (σ) maps to s.

To handle the problem of ambiguity, Transliterator uses an apostrophe to disambiguate character sequences. Using this procedure, the Greek character PI SIGMA ($\pi\sigma$) maps to p's. In Japanese, whenever an ambiguous sequence in the target script does not result from a single letter, the transform uses an apostrophe to disambiguate it. For example, it uses this procedure to distinguish between man'ichi and manichi.

Note: Some characters in a target script are not normally found outside of certain contexts. For example, the small Japanese "ya" character, as in "kya" (キヤ), is not normally found in isolation. To handle such characters, Transliterator uses a tilde. For example, the input "~ya" would produce an isolated small "ya". When transliterating to Greek, the input "a~s" would produce a non-final Greek sigma ($\alpha\sigma$) at the end of a word. Likewise, the input "~sa" would produce a final sigma in a non-final position ($\zeta\alpha$).

For the general script transforms, a common technique for reversibility is to use extra accents to distinguish between letters that may not be otherwise distinguished. For example, the following shows Greek text that is mapped to fully reversible Latin:

Input

Field Name	Description										
Any string field	The Transliterator stage can transliterate any string field. You can specify which fields to transliterate in the Transliterator stage options.										
TransliteratorID	<p>Overrides the default transliteration specified in the Transliterator stage options. Use this field if you want to specify a different transliteration for each record.</p> <p>For Example:</p> <table border="0"> <tr> <td>Arabic-Latin</td> <td>From Arabic to Latin.</td> </tr> <tr> <td>Greek-Latin</td> <td>From Greek to Latin.</td> </tr> <tr> <td>Latin-Hangul</td> <td>From Latin to Hangul.</td> </tr> <tr> <td>Latin-Katakana</td> <td>From Latin to Katakana.</td> </tr> <tr> <td>Fullwidth-Halfwidth</td> <td>From full width to half width.</td> </tr> </table>	Arabic-Latin	From Arabic to Latin.	Greek-Latin	From Greek to Latin.	Latin-Hangul	From Latin to Hangul.	Latin-Katakana	From Latin to Katakana.	Fullwidth-Halfwidth	From full width to half width.
Arabic-Latin	From Arabic to Latin.										
Greek-Latin	From Greek to Latin.										
Latin-Hangul	From Latin to Hangul.										
Latin-Katakana	From Latin to Katakana.										
Fullwidth-Halfwidth	From full width to half width.										

Options

Table 27: Transliterator Options

Option	Description/Valid Values
From	<p>The script used by the fields that you want to transliterate. For a description of the supported scripts, see Transliterator on page 338.</p> <p>Note: The Transliterator stage does not support transliteration between all scripts. The From and To fields automatically reflect the valid values based on your selection.</p>
To	<p>The script that you want to convert the field into. For a description of the supported scripts, see Transliterator on page 338.</p> <p>Note: The Transliterator stage does not support transliteration between all scripts. The From and To fields automatically reflect the valid values based on your selection.</p>
Swap button	<p>Click the swap button to exchange the languages in the From and To fields.</p>
	
Fields to transliterate	<p>Specifies the fields that you want to transliterate.</p>
Remove Accent Marks	<p>Select the check box to remove accent marks from the field. This option remains turned off by default.</p>

Output

The Transliterator stage transliterates the fields you specify. It does not produce any other output.

Universal Name

Spectrum Universal Name

To perform the most accurate standardization you may need to break up strings of data into multiple fields. Spectrum Technology Platform provides advanced parsing features that enable you to parse personal names, company names, and many other terms and abbreviations. In addition, you can create your own list of custom terms to use as the basis of scan and extract operations.

Components

Spectrum Universal Name consists of:

- **Name Variant Finder**—This component works in either a first name or last name mode in order to query the names database and return name variations. Name Variant Finder can limit the number of results by gender and culture. Spectrum Technology Platform includes a base names file. Add-on names files of other cultures are available and can be deployed by copying the appropriate JAR file into the `modules/tables/ext` folder. Arabic names are available as an add-on. The data in this file is by Nomino, Inc.
- **Open Name Parser**—The open name parser breaks down personal and business names and other terms in the name data field into their component parts. These parsed name elements are then subsequently available to other automated operations such as name matching, name standardization, or multi-record name consolidation.

Name Parser (DEPRECATED)

Attention: The Name Parser stage is deprecated and may not be supported in future releases. Use Open Name Parser for parsing names.

Name Parser breaks down personal and business names and other terms in the name data field into their component parts. The parsing process includes an explanation of the function, form and syntactical relationship of each part to the whole. These parsed name elements are then subsequently available to other automated operations such as name matching, name standardization, or multirecord name consolidation.

Name parsing does the following:

- Determines the entity type of a name in order to describe the function which the name performs. Name entity types are divided into two major groupings: Personal names and business names with subgroups within these major groupings.
- Determines the form of a name in order to understand which syntax the parser should follow for parsing. Personal names usually take on a natural (signature) order or a reverse order. Business names are usually ordered hierarchically.
- Determines and labels the component parts of a name so that the syntactical relationship of each name part to the entire name is identified. The personal name syntax includes prefixes, first, middle, and last name parts, suffixes and account description terms among other personal name parts. The business name syntax includes the primary text, insignificant terms, prepositions, objects of the preposition and suffix terms among other business name parts.
- Determines the gender of the name. The gender is determined based on cultural assumptions which you specify. For example, Jean is a male name in France but a female name in the U.S. If you know the names you are processing are from France, you could specify French as the gender determination culture. The Name Parser uses data from the First Name and Compound First Names tables to determine gender. If a name is not found in either table and a title is present in the name, the parser checks the Title table to determine gender. Otherwise, the gender is marked as unknown.

Note: If a field on your input record already contains one of the supported cultures, you can predefine the GenderDeterminationSource field in your input to override the Gender Determination Source in the GUI.

- Assigns a parsing score which indicates the degree of confidence which the parser has that its parsing is correct.

Input

Attention: The Name Parser stage is deprecated and may not be supported in future releases. Use Open Name Parser for parsing names.

Table 28: Name Parser Input

Field Name	Description / Valid Values																										
GenderDeterminationSource	<p>The culture of the name data to use to determine gender. Default uses cross-cultural rules. For example, Jean is commonly a female name and Default identifies it as such, but it is identified as a male name if you select French. The options are listed below along with example countries for each culture. Note that the list of countries under each culture is not exhaustive.</p> <table border="0"> <tr> <td>SLAVIC</td> <td>Bosnia, Poland, Albania.</td> </tr> <tr> <td>ARMENIAN</td> <td>Armenia.</td> </tr> <tr> <td>DEFAULT</td> <td>Bulgaria, Cayman Islands, Ireland, U.S., U.K.</td> </tr> <tr> <td>FRENCH</td> <td>France.</td> </tr> <tr> <td>SCANDINAVIAN</td> <td>Denmark, Finland, Iceland, Norway, Sweden.</td> </tr> <tr> <td>GERMANIC</td> <td>Austria, Germany, Luxembourg, Switzerland, The Netherlands.</td> </tr> <tr> <td>GREEK</td> <td>Greece.</td> </tr> <tr> <td>HUNGARIAN</td> <td>Hungary.</td> </tr> <tr> <td>ITALIAN</td> <td>Italy.</td> </tr> <tr> <td>PORTUGUESE</td> <td>Portugal.</td> </tr> <tr> <td>ROMANIA</td> <td>Romania.</td> </tr> <tr> <td>HISPANIC</td> <td>Spain.</td> </tr> <tr> <td>ARABIC</td> <td>Tunisia.</td> </tr> </table> <p>GenderDeterminationSource is also used by Name Variant Finder to limit the returned name variations based on culture. For more information, see Name Variant Finder on page 366.</p>	SLAVIC	Bosnia, Poland, Albania.	ARMENIAN	Armenia.	DEFAULT	Bulgaria, Cayman Islands, Ireland, U.S., U.K.	FRENCH	France.	SCANDINAVIAN	Denmark, Finland, Iceland, Norway, Sweden.	GERMANIC	Austria, Germany, Luxembourg, Switzerland, The Netherlands.	GREEK	Greece.	HUNGARIAN	Hungary.	ITALIAN	Italy.	PORTUGUESE	Portugal.	ROMANIA	Romania.	HISPANIC	Spain.	ARABIC	Tunisia.
SLAVIC	Bosnia, Poland, Albania.																										
ARMENIAN	Armenia.																										
DEFAULT	Bulgaria, Cayman Islands, Ireland, U.S., U.K.																										
FRENCH	France.																										
SCANDINAVIAN	Denmark, Finland, Iceland, Norway, Sweden.																										
GERMANIC	Austria, Germany, Luxembourg, Switzerland, The Netherlands.																										
GREEK	Greece.																										
HUNGARIAN	Hungary.																										
ITALIAN	Italy.																										
PORTUGUESE	Portugal.																										
ROMANIA	Romania.																										
HISPANIC	Spain.																										
ARABIC	Tunisia.																										
Name	The name you want to parse. This field is required.																										

Options

Attention: The Name Parser stage is deprecated and may not be supported in future releases. Use Open Name Parser for parsing names.

To specify the Name Parser options, double-click the instance of Name Parser on the canvas. The Name Parser Options dialog displays.

Table 29: Name Parser Options

Option	Description
Parse personal names	Check this box to parse personal names.
Separate conjoined names into multiple records Select a match results in the Match Results List and then click Remove .	Click this box to separate names containing more than one individual into multiple records, for example, <i>Bill & Sally Smith</i> . When a conjoined record results in two separate name records, a Parser Record ID output field is generated. Each pair of separate name records are identified with the same Parser Record ID.
Gender Determination Source Select a match results in the Match Results List and then click Remove .	Determines how the Name Parser assigns a gender to the name. For most cases, Default is the best setting because it covers a wide variety of names. If you are processing names from a specific culture, select that culture. Selecting a specific culture helps ensure that the proper gender is assigned to the names. For example, if you leave Default selected, then the name Jean is identified as a female name. If you select French, it is identified as a male name. Note: If you select a culture but the name is not found in that culture, gender is determined using the Default culture, which includes data from a variety of cultures.
Order	Specifies how the name fields are ordered in your input records. One of the following: Natural The name fields are ordered by Title, First Name, Middle Name, Last Name, and Suffix. Reverse The name fields are ordered by Last Name first. Mixed The name fields are ordered using a combination of natural and reverse.
Retain Periods	Retains punctuation in the parsed personal name field.
Parse Business Names	Check this box to parse business names.
Retain Periods	Check this box to return punctuation to the parsed business name field.

Option	Description
User-Defined Table	Click any of the User-Defined Tables to add values to existing values in the various parser tables. This capability enables you to customize tables for your unique business environment. Click Configure to select an XML file that contains the values that you want to add. For more information about user-defined tables, see Modifying Name Parser User-Defined Tables on page 348.

Modifying Name Parser User-Defined Tables

Attention: The Name Parser stage is deprecated and may not be supported in future releases. Use Open Name Parser for parsing names.

You can add, modify, and delete values in the Name Parser tables to customize them for your unique business environment.

Name Parser's user-defined tables are XML files located by default in the <Drive>:\Program Files\Precisely\Spectrum\server\modules\parser\data folder. Spectrum Technology Platform includes the following user-defined tables:

UserAccountDescriptions.xml

Table 30: UserAccountDescriptions.xml Columns

Column Name	Description / Valid Values
LookupValue	A lookup term commonly found in an Account Description. Any single-word text. Case insensitive.

Example entry:

```
<table-data>
  <deleted-entries delimiter-character="|">
    <deleted-entry-group>
      <![CDATA[
        LookupValue
        ART
        AND
      ]]>
    </deleted-entry-group>
  </deleted-entries>
```

```
<added-entries delimiter-character="|">
  <![CDATA[
    LookupValue
    A/C
    ACCOUNT
    EXP
  ]]>
</added-entries>
</table-data>
```

UserCompanyPrepositions.xml

Table 31: UserCompanyPrepositions.xml Columns

Column Name	Description / Valid Values
LookupValue	Any preposition (for example, "of" or "on") commonly found in company names. Any single-word text. Case insensitive.

Example entry:

```
<table-data>
  <deleted-entries delimiter-character="|">
    <deleted-entry-group>
      <![CDATA[
        LookupValue
        AROUND
        NEAR
      ]]>
    </deleted-entry-group>
  </deleted-entries>
  <added-entries delimiter-character="|">
    <![CDATA[
      LookupValue
      ABOUT
      AFTER
      ACROSS
    ]]>
  </added-entries>
</table-data>
```

*UserCompanySuffixes.xml***Table 32: UserCompanySuffixes.xml Columns**

Column Name	Description / Valid Values
LookupValue	Any suffix commonly found in company names. Examples include "Inc." and "Co." Any single-word text. Case insensitive.

Example entry:

```
<table-data>
  <deleted-entries delimiter-character="|">
    <deleted-entry-group>
      <![CDATA[
        LookupValue
        SANDY
        CLUE
      ]]>
    </deleted-entry-group>
  </deleted-entries>
  <added-entries delimiter-character="|">
    <![CDATA[
      LookupValue
      LTD
      LLC
      CO
      INC
    ]]>
  </added-entries>
</table-data>
```

*UserCompanyTerms.xml***Table 33: UserCompanyTerms.xml Columns**

Column Name	Description / Valid Values
LookupValue	Any term commonly found in a company name. Any single-word text. Case insensitive.

Example entry:

```
<table-data>
  <deleted-entries delimiter-character="|">
    <deleted-entry-group>
      <![CDATA[
        LookupValue
        MARY
        BLUE
      ]]>
    </deleted-entry-group>
  </deleted-entries>
  <added-entries delimiter-character="|">
    <![CDATA[
      LookupValue
      ARC
      ARCADE
      ASSEMBLY
      ARIZONA
    ]]>
  </added-entries>
</table-data>
```

UserCompoundFirstNames.xml

This table contains user-defined compound first names. Compound names are names that consist of two words.

Table 34: UserCompoundFirstNames.xml Columns

Column Name	Description / Valid Values
FirstName	The compound first name. Maximum of two words. Case insensitive.
Culture	The culture in which this FirstName/Gender combination applies. You may use any of the values that are valid in the GenderDeterminationSource input field. For more information, see Input on page 345.

Column Name	Description / Valid Values
Gender	<p>The gender most commonly associated with this FirstName/Culture combination. One of the following:</p> <p>M The name is a male name.</p> <p>F The name is a female name.</p> <p>A Ambiguous. The name can be either male or female.</p> <p>U Unknown. The gender of this name is not known. Unknown is assumed if this field is left blank.</p>
Frequency	Not used in this release. You may leave this column blank.

Example entry:

```

<table-data>
  <deleted-entries delimiter-character="|">
    <deleted-entry-group>
      <![CDATA[
        FirstName
        ANN MARIE
        BILLY JOE
      ]]>
    </deleted-entry-group>
    <deleted-entry-group>
      <![CDATA[
        FirstName|Frequency
        KAREN SUE|0.126
        BILLY JOE|0.421
      ]]>
    </deleted-entry-group>
    <deleted-entry-group>
      <![CDATA[
        FirstName|Gender|Culture
        JEAN ANN|M|DEFAULT
        JEAN CLUADE|F|FRENCH
      ]]>
    </deleted-entry-group>
  </deleted-entries>
  <added-entries delimiter-character="|">
    <![CDATA[
      FirstName|Gender|Culture
      JOHN Henry|M|DEFAULT
      A'SHA A'MAR|F|ARABIC
      BILLY JO|A|DEFAULT
    ]]>
  </added-entries>

```

```

]]>
</added-entries>
</table-data>

```

UserConjunctions.xml

This table contains a list of user-defined conjunctions, such as "and", "or", or "&".

Table 35: UserConjunctions.xml Columns

Column Name	Description / Valid Values
LookupValue	Any conjunction. Must be a single word. Case insensitive.

Example entries:

```

<table-data>
  <deleted-entries delimiter-character="|">
    <deleted-entry-group>
      <![CDATA[
        LookupValue
        FIND
        CARE
        %
      ]]>
    </deleted-entry-group>
  </deleted-entries>
  <added-entries delimiter-character="|">
    <![CDATA[
      LookupValue
      &
      AND
      OR
    ]]>
  </added-entries>
</table-data>

```

UserFirstNames.xml

Table 36: UserFirstNames.xml Columns

Column Name	Description / Valid Values
FirstName	The first name described by this table row. Case insensitive.
Gender	<p>The gender most commonly associated with this FirstName/Culture combination. One of the following:</p> <p>M The name is a male name.</p> <p>F The name is a female name.</p> <p>A Ambiguous. The name can be either male or female.</p> <p>U Unknown. The gender of this name is not known. Unknown is assumed if this field is left blank.</p>
Culture	The culture in which this FirstName/Gender combination applies. You may use any of the values that are valid in the GenderDeterminationSource input field. For more information, see Input on page 345.

Example entry:

```

<table-data>
  <deleted-entries delimiter-character="|">
    <deleted-entry-group>
      <![CDATA[
        FirstName
        AADEL
        AADIL
      ]]>
    </deleted-entry-group>
    <deleted-entry-group>
      <![CDATA[
        FirstName
        A ' SACE
        A ' BOCKETT
      ]]>
    </deleted-entry-group>
    <deleted-entry-group>
      <![CDATA[
        FirstName|Gender|Culture
      ]]>
    </deleted-entry-group>
  </deleted-entries>
</table-data>

```

```

                ALII|M|DEFAULT
                AISHA|F|ARABIC
            ]]>
        </deleted-entry-group>
    </deleted-entry-group>
    <![CDATA[
        FirstName|Gender
        JOHE|M
    ]]>
    </deleted-entry-group>
</deleted-entries>
<added-entries delimiter-character="|">
    <![CDATA[
        FirstName|Gender|Culture
        JOHE|M|DEFAULT
        A'SHAN|F|ARABIC
    ]]>
</added-entries>
</table-data>

```

UserGeneralSuffixes.xml

This table contains a list of user-defined suffixes used in personal names that are not maturity suffixes, such as "MD" or "PhD".

Table 37: UserGeneralSuffixes.xml Columns

Column Name	Description / Valid Values
LookupValue	Any suffix that is frequently applied to personal names and is not a maturity suffix. Must be a single word. Case insensitive.

Example entry:

```

<table-data>
  <deleted-entries delimiter-character="|">
    <deleted-entry-group>
      <![CDATA[
        LookupValue
        AND
        WILL
        TUNA
      ]]>
    </deleted-entry-group>

```

```

</deleted-entries>
<added-entries delimiter-character="|" ">
  <![CDATA[
    LookupValue
    ACCOUNTANT
    ATTORNEY
    ANALYST
    ASSISTANT
  ]]>
</added-entries>
</table-data>

```

UserLastNamePrefixes.xml

This table contains a list of user-defined prefixes that occur in a person's last name such as "Van", "De", or "La".

Table 38: UserLastNamePrefixes.xml Columns

Column Name	Description / Valid Values
LookupValue	Any prefix that occurs as part of an individual's last name. Any single-word text. Case insensitive.

Example entry:

```

<table-data>
  <deleted-entries delimiter-character="|" ">
    <deleted-entry-group>
      <![CDATA[
        LookupValue
        DO
        RUN
        ANIMAL
      ]]>
    </deleted-entry-group>
  </deleted-entries>
  <added-entries delimiter-character="|" ">
    <![CDATA[
      LookupValue
      D'
      DA
      DEN
      DEL
    ]]>
  </added-entries>
</table-data>

```

```

]]>
</added-entries>
</table-data>

```

UserLastNames.xml

Table 39: UserLastNames.xml Columns

Column Name	Description / Valid Values
LastName	The last name described by this table row. Case insensitive.
Gender	<p>The gender most commonly associated with this FirstName/Culture combination. One of the following:</p> <p>M The name is a male name.</p> <p>F The name is a female name.</p> <p>A Ambiguous. The name can be either male or female.</p> <p>U Unknown. The gender of this name is not known. Unknown is assumed if this field is left blank.</p>
Culture	The culture in which this FirstName/Gender combination applies. You may use any of the values that are valid in the GenderDeterminationSource input field. For more information, see Input on page 345.

Example entry:

```

<table-data>
  <deleted-entries delimiter-character="|">
    <deleted-entry-group>
      <![CDATA[
        LastName
        Rusod
        AADIL
      ]]>
    </deleted-entry-group>
    <deleted-entry-group>
      <![CDATA[
        LastName

```

```

                KAASEEY
                JOIEN
            ]]>
        </deleted-entry-group>
    </deleted-entries>
    <added-entries delimiter-character="|">
        <![CDATA[
            LastName|Culture|Gender
            SMITH|ENGLISH|A
            WILSON|ENGLISH|A
            JONES|ENGLISH|A
        ]]>
    </added-entries>
</table-data>

```

UserMaturitySuffixes.xml

This table contains user-defined generational suffixes used in a person's name, such as "Jr." or "Sr."

Table 40: UserMaturitySuffixes.xml Columns

Column Name	Description / Valid Values
LookupValue	A generational suffix used in personal names. Any single-word text. Case insensitive.

Example entry:

```

<table-data>
    <deleted-entries delimiter-character="|">
        <deleted-entry-group>
            <![CDATA[
                LookupValue
                I
                V
                18
                VI
            ]]>
        </deleted-entry-group>
    </deleted-entries>
    <added-entries delimiter-character="|">
        <![CDATA[
            LookupValue
            I
            II
        ]]>
    </added-entries>
</table-data>

```

```

        III
      ]]>
    </added-entries>
  </table-data>

```

UserTitles.xml

This table contains user-defined titles used in a person's name, such as "Mr." or "Ms."

Table 41: UserTitles.xml Columns

Column Name	Description / Valid Values
LookupValue	A title used in personal names. Any single-word text. Case insensitive.
Gender	The gender most commonly associated with this title. One of the following: <ul style="list-style-type: none"> M The name is a male name. F The name is a female name. A Ambiguous. The name can be either male or female. U Unknown. The gender of this name is not known. Unknown is assumed if this field is left blank.

Example entry:

```

<table-data>
  <deleted-entries delimiter-character="|">
    <deleted-entry-group>
      <![CDATA[
        LookupValue
        Belt
        Friend
        Thursday
        Red
      ]]>
    </deleted-entry-group>
  </deleted-entries>
  <added-entries delimiter-character="|">
    <![CDATA[
      LookupValue|Gender
      Mrs|F
    ]]>
  </added-entries>
</table-data>

```

```

        Mr|M
        Most|F
    ]]>
</added-entries>
</table-data>

```

Sample User-Defined Table

The figure below shows a sample UserFirstNames.xml table and the syntax to use when modifying user-defined tables.

```

<table-data>
  <deleted-entries delimiter-character="|">
    <deleted-entry-group>
      <![CDATA[
        FirstName
        AADEL
        AADIL
      ]]>
    </deleted-entry-group>
    <deleted-entry-group>
      <![CDATA[
        FirstName|Frequency
        A'SACE|0.126
        A'BECKETT|0.421
      ]]>
    </deleted-entry-group>
    <deleted-entry-group>
      <![CDATA[
        FirstName|Gender|Culture|VariantGroup
        ALI|M|DEFAULT|GROUP88
        AISHA|F|ARABIC|GROUP43
      ]]>
    </deleted-entry-group>
    <deleted-entry-group>
      <![CDATA[
        FirstName|Gender
        JOHN|M
      ]]>
    </deleted-entry-group>
  </deleted-entries>
  <added-entries delimiter-character="|">
    <![CDATA[
      FirstName|Gender|Culture
      JOHN|M|DEFAULT
      A'SHA|F|ARABIC
      JAMES|M|DEFAULT
    ]]>
  </added-entries>

```

</table-data>

Output

Attention: The Name Parser stage is deprecated and may not be supported in future releases. Use Open Name Parser for parsing names.

Table 42: Name Parser Output

Field Name	Format	Description / Valid Values
AccountDescription	String	An account description that is part of the name. For example, in "Mary Jones Account # 12345", the account description is "Account#12345".
EntityType	String	Indicates the type of name. One of the following: Firm The name is a company name. Personal The name is an individual person's name.
Fields Related to Names of Companies		
FirmModifier.1.Object	String	The first object of a preposition occurring in firm name. For example, in the firm name "Pratt & Whitney Division of United Technologies", the first object of a preposition is "United Technologies".
FirmModifier.1.Preposition	String	The first preposition occurring in firm name. For example, in the firm name "Pratt & Whitney Division of United Technologies", "of" would be the first preposition.
FirmModifier.2.Object	String	The second object of a preposition occurring in firm name. For example, in the firm name "Church of Our Lady of Lourdes", the second object of a preposition is the second "Lourdes".
FirmModifier.2.Preposition	String	The second preposition occurring in firm name. For example, in the firm name "Church of Our Lady of Lourdes", the second preposition is the second "of".

Field Name	Format	Description / Valid Values
FirmName	String	The name of a company. For example, "Precisely Software, Inc."
FirmPrimary	String	The base part of a company's name. For example, "Precisely Software".
FirmSuffix	String	The corporate suffix. For example, "Co." and "Inc."
Fields Related to Names of Individual People		
FirstName	String	The first name of a person.
FirstNameVariantGroup	String	<p>A numeric ID that indicates the group of similar names to which first name belongs. For example, Muhammad, Mohammed, and Mehmet all belong to the same Name Variant Group. The actual group ID is assigned when the add-on data is loaded.</p> <p>This field is only populated if you have purchased the Name Variant Group feature.</p>
GenderCode	String	<p>A person's gender as determined by analyzing the first name. One of the following:</p> <ul style="list-style-type: none"> A Ambiguous. The name is both a male and a female name. For example, Pat. F Female. The name is a female name. M Male. The name is a male name. U Unknown. The name could not be found in the gender table.
GenderDeterminationSource	String	The culture used to determine a name's gender. If the name could not be found in the gender table, this field is blank.
GeneralSuffix	String	A person's general/professional suffix. For example, MD or PhD.
LastName	String	The last name of a person.

Field Name	Format	Description / Valid Values
MaturitySuffix	String	A person's maturity/generational suffix. For example, Jr. or Sr.
MiddleName	String	The middle name of a person.
NameScore	String	Score representing quality of the parsing operation, from 0 to 100. 0 indicates poor quality and 100 indicates high quality.
ParserRecordID	String	A unique ID assigned to each input record.
TitleOfRespect	String	A person's title, such as Mr., Mrs., Dr., or Rev.
Fields Related to Conjoined Names		
PersonalName.2.FirstName	String	The first name of the second person in a conjoined name. An example of a conjoined name is "John and Jane Smith."
PersonalName.2.FirstNameVariantGroup	String	<p>A numeric ID that indicates the group of similar names to which first name of the second person in a conjoined name belongs. For example, Muhammad, Mohammed, and Mehmet all belong to the same Name Variant Group. The actual group ID is assigned when the add-on data is loaded.</p> <p>This field is only populated if you have purchased the Name Variant Group feature.</p>
PersonalName.2.GenderCode	String	<p>The gender of the second person in a conjoined name as determined by Name Parser analyzing the first name. An example of a conjoined name is "John and Jane Smith." One of the following:</p> <p>A Ambiguous. The name is both a male and a female name. For example, Pat.</p> <p>F Female. The name is a female name.</p> <p>M Male. The name is a male name.</p> <p>U Unknown. The name could not be found in the gender table.</p>

Field Name	Format	Description / Valid Values
PersonalName2GenderDeterminationSource	String	The culture used to determine the gender of the second person in a conjoined name. An example of a conjoined name is "John and Jane Smith."
PersonalName.2.GeneralSuffix	String	The general/professional suffix of the second person in a conjoined name. An example of a conjoined name is "John and Jane Smith." Examples of general suffixes are MD and PhD.
PersonalName.2.LastName	String	The last name of the second person in a conjoined name. An example of a conjoined name is "John and Jane Smith."
PersonalName.2.MaturitySuffix	String	The maturity/generational suffix of the second person in a conjoined name. An example of a conjoined name is "John and Jane Smith." Examples of maturity suffixes are Jr. and Sr.
PersonalName.2.MiddleName	String	The middle name of the second person in a conjoined name. An example of a conjoined name is "John and Jane Smith."
PersonalName.2.TitleOfRespect	String	The title of respect for the second name in a conjoined name. For example, "Mr. and Mrs. Smith" is a conjoined name. Examples of titles of respect are Mr., Mrs., and Dr.
PersonalName.3.FirstName	String	The first name of the third person in a conjoined name. For example, "Mr. & Mrs. John Smith & Dr. Mary Jones" is a conjoined name.
PersonalName.3.FirstNameVariantGroup	String	<p>A numeric ID that indicates the group of similar names to which first name of the second person in a conjoined name belongs. For example, Muhammad, Mohammed, and Mehmet all belong to the same Name Variant Group. The actual group ID is assigned when the add-on data is loaded.</p> <p>This field is only populated if you have purchased the Name Variant Group feature.</p>

Field Name	Format	Description / Valid Values
PersonalName.3.GenderCode	String	<p>The gender of the third person in a conjoined name as determined by Name Parser analyzing the first name. An example of a conjoined name is "Mr. & Mrs. John Smith & Adam Jones". One of the following:</p> <p>A Ambiguous. The name is both a male and a female name. For example, Pat.</p> <p>F Female. The name is a female name.</p> <p>M Male. The name is a male name.</p> <p>U Unknown. The name could not be found in the gender table.</p>
PersonalName.3.GenderDeterminationSource	String	The culture used to determine the gender of the third person in a conjoined name. "Mr. & Mrs. John Smith & Adam Jones".
PersonalName.3.GeneralSuffix	String	The general/professional suffix of the third person in a conjoined name. An example of a conjoined name is "Mr. & Mrs. John Smith & Adam Jones PhD." Examples of general suffixes are MD and PhD.
PersonalName.3.LastName	String	The last name for the third person in a conjoined name. For example, "Mr. & Mrs. John Smith & Dr. Mary Jones" is a conjoined name.
PersonalName.3.MaturitySuffix	String	The maturity/generational suffix of the third person in a conjoined name. An example of a conjoined name is "Mr. & Mrs. John Smith & Adam Jones Sr." Examples of maturity suffixes are Jr. and Sr.
PersonalName.3.MiddleName	String	The middle name for the third person in a conjoined name. For example, "Mr. & Mrs. John Smith & Dr. Mary Jones" is a conjoined name.
PersonalName.3.TitleOfRespect	String	The title of respect for the third name in a conjoined name. For example, "Mr. & Mrs. John Smith & Dr. Mary Jones" is a conjoined name. Examples of titles of respect are Mr., Mrs., and Dr.

Name Variant Finder

Name Variant Finder works in either first name or last name mode to query a database to return alternative versions of a name. For example, "John" and "Jon" are variants for the name "Johnathan". Name Variant Finder requires add-on dictionaries that can be installed using the Universal Name, Data Normalization, and Advanced Matching database load utility. Contact your sales representative for information on how to obtain these optional culture-specific dictionaries.

Input

Table 43: Name Variant Finder Input Fields

Field Name	Description / Valid Values								
FirstName	The name for which you want to find variants, if the name is a given name.								
LastName	The name for which you want to find variants, if the name is a surname.								
GenderCode	<p>The gender of the name in the FirstName field. One of the following:</p> <p>Note: Gender codes only apply to first names, not last names.</p> <table border="0"> <tr> <td>M</td> <td>The name is a male name.</td> </tr> <tr> <td>F</td> <td>The name is a female name.</td> </tr> <tr> <td>A</td> <td>Ambiguous. The name can be either male or female.</td> </tr> <tr> <td>U</td> <td>Unknown. The gender of this name is not known.</td> </tr> </table>	M	The name is a male name.	F	The name is a female name.	A	Ambiguous. The name can be either male or female.	U	Unknown. The gender of this name is not known.
M	The name is a male name.								
F	The name is a female name.								
A	Ambiguous. The name can be either male or female.								
U	Unknown. The gender of this name is not known.								
Ethnicity	<p>The culture most commonly associated with the name in the FirstName or LastName field. You can use the Name Parser or Open Parser stages to populate this field if you do not know the ethnicity for a name.</p> <p>Note: This field was formerly named GenderDeterminationSource.</p>								

Options

Table 44: Name Variant Finder Options

Option	Description
First Name	Finds name variations based on first name.
Last Name	Finds name variations based on last name.
Gender Code	Returns the name variations only for the gender specified in the record's GenderCode field. For information about the GenderCode field, see Input on page 366.
Ethnicity	Returns name variations only for the culture specified in the record's Ethnicity field. For information about the Ethnicity field, see Input on page 366.
Romanized	Returns the English romanized version of the name. A romanized name is one that has been converted from a non-Latin script to the Latin script. For example, Achin is the Romanized version of the Korean name 아진.
Native	Returns the name in the native script of the name's culture. For example, a Korean name would be returned in Hangul.
Kana	If you select Native , you can choose to return Japanese names in Kana by selecting this option. Kana is comprised of hiragana and katakana scripts. Note: You must have licensed the Asian Plus Pack database to look up Japanese name variants. For more information, contact your sales executive.
Kanji	If you select Native , you can choose to return Japanese names in Kanji by selecting this option. Kanji is one of the scripts used in the Japanese language. Note: You must have licensed the Asian Plus Pack database to look up Japanese name variants. For more information, contact your sales executive.

Output

Table 45: Name Variant Finder Outputs

Field Name	Format	Description / Valid Values
CandidateGroup	String	Identifies a grouping of an input name and its name variations. Each input name is given a CandidateGroup number. The variations for that input name are given the same CandidateGroup number.
Ethnicity	String	The culture of a name determined by the Core Name and add-on dictionaries. Note: This field was formerly named GenderDeterminationSource.
FirstName	String	The given name of a person.
GenderCode	String	The gender of a name determined by the Core Name and add-on dictionaries. One of the following: M The name is a male name. F The name is a female name. A Ambiguous. The name can be either male or female. U Unknown. The gender of this name is not known.
LastName	String	The last name of a person.
TransactionalRecordType	String	Specifies how the name was used in the matching process. One of the following: Suspect A suspect record is used as input to a query. Candidate A candidate record is a result returned from a query.

Open Name Parser

Open Name Parser breaks down personal and business names and other terms in the name data field into their component parts. These parsed name elements are then subsequently available to other automated operations such as name matching, name standardization, or multiple-record name consolidation.

Open Name Parser does the following:

- Determines the type of a name in order to describe the function that the name performs. Name entity types are divided into two major groups: personal names and business names. Within each of these major groups are subgroups.
- Determines the form of a name in order to understand which syntax the parser should follow for parsing. Personal names usually take on a natural (signature) order or a reverse order. Business names are usually ordered hierarchically.
- Determines and labels the component parts of a name so that the syntactical relationship of each name part to the entire name is identified. The personal name syntax includes prefixes, first, middle, and last name parts, suffixes, and account description terms, among other personal name parts. The business name syntax includes the firm name and suffix terms.
- Parses conjoined personal and business names and either retains them as one record or splits them into multiple records. Examples of conjoined names include "Mr. and Mrs. John Smith" and "Baltimore Gas & Electric dba Constellation Energy".
- Parses output as records or as a list.
- Enables you to use the Open Parser Domain Editor to create new domains that can be used in the Open Name Parser Advanced Options.
- Assigns a parsing score that reflects the degree of confidence that the parsing is correct.

Input

Table 46: Open Name Parser Input

Field Name	Description								
CultureCode	<p>The culture of the input name data. The options are listed below.</p> <table border="0"> <tr> <td>Null (empty)</td> <td>Global culture (default).</td> </tr> <tr> <td>de</td> <td>German.</td> </tr> <tr> <td>es</td> <td>Spanish.</td> </tr> <tr> <td>ja</td> <td>Japanese.</td> </tr> </table> <p>Note: If you added your own domain using the Open Parser Domain Editor, the cultures and culture codes for that domain are also valid.</p>	Null (empty)	Global culture (default).	de	German.	es	Spanish.	ja	Japanese.
Null (empty)	Global culture (default).								
de	German.								
es	Spanish.								
ja	Japanese.								
Name	The name you want to parse. This field is required.								

Options

Open Name Parser options can be configured at the stage level, through any of the Spectrum Technology Platform clients, or at runtime, using dataflow options.

Parsing Options

The following table lists the options that control the parsing of names.

Table 47: Open Name Parser Parsing Options

Option Name	Description
Parse personal names	<p>Specifies whether to parse personal names.</p> <p>Natural The name fields are ordered by Title, First Name, Middle Name, Last Name, and Suffix.</p> <p>Reverse The name fields are ordered by Last Name first.</p> <p>Both The name fields are ordered using a combination of natural and reverse.</p>
Conjoined names	Specifies whether to parse conjoined names.
Split conjoined names into multiple records	<p>Specifies whether to separate names containing more than one individual into multiple records, for example, Bill & Sally Smith.</p> <p>Use a Unique ID Generator stage to create an ID for each of the split records.</p>
Parse business names	Specifies whether to parse business names.
Output results as list	Specifies whether to return the parsed name elements in a list form.
Shortcut threshold	<p>Specifies how to balance performance versus quality. A faster performance will result in lower quality output; likewise, higher quality will result in slower performance. When this threshold is met, no other processing will be performed on the record.</p> <p>The default is 100.</p>

Cultures Options

The following table lists the options that control name cultures.

Table 48: Open Name Parser Cultures Options

Option Name	Description
Cultures	<p>Specifies which culture(s) you want to include in the parsing grammar. Global Culture is the default selection.</p> <p>Note: If you added your own domain using the Open Parser Domain Editor, the cultures and culture codes for that domain will appear here as well.</p> <p>Click the Up and Down buttons to set the order in which you want the cultures to run.</p>

Advanced Options

The following table lists the advanced options for name parsing.

Table 49: Open Name Parser Advanced Options

Option	Description
Advanced Options	<p>Use the Domain drop-down to select the appropriate domain for each Name.</p> <p>Click the Up and Down buttons to set the order in which you want the parsers to run. Results will be returned for the first domain that scores higher than the number set in the Shortcut threshold field. If no domain reaches that threshold, results for the domain with the highest score are returned. If multiple domains reach the threshold at the same time, priority goes to the domain that was run first (determined by the order set here) and its results will be returned.</p> <p>Note: If you added your own domain using the Open Parser Domain Editor, that domain will appear here as well.</p>

Configuring Options at Runtime

Open Name Parser options can be configured and passed at runtime if they are exposed as dataflow options. This enables you to override the existing configuration with JSON-formatted name parsing

strings. You can also set stage options when calling the job through a process flow or through the job executor command-line tool.

To define Open Name Parser options at runtime:

1. In Enterprise Designer, open a dataflow that uses the Open Name Parser stage.
2. Save and expose that dataflow.
3. Go to `Edit > Dataflow Options`.
4. In the **Map dataflow options to stages** table, expand Open Name Parser and edit options as necessary. Check the box for the option you want to edit, then change the value in the **Default value** drop-down.
5. Optional: Change the name of the options in the **Option label** field.
6. Click **OK** twice.

Output

Table 50: Open Name Parser Output

Field Name	Format	Description
AccountDescription	String	An account description that is part of the name. For example, in "Mary Jones Account # 12345", the account description is "Account#12345".
Names	String	A hierarchical field that contains a list of parsed elements. This field is returned when you check the Output results as list box under Parsing Options.
Fields Related to Names of Companies		
FirmConjunction	String	Indicates that the name of a firm contains a conjunction such as "d/b/a" (doing business as), "o/a" (operating as), and "t/a" (trading as).
FirmName	String	The name of a company. For example, <code>Precisely</code> .
FirmSuffix	String	The corporate suffix. For example, "Co." and "Inc."
IsFirm	String	Indicates that the name is a firm rather than an individual.

Field Name	Format	Description
Fields Related to Names of Individual People		
Conjunction	String	Indicates that the name contains a conjunction such as "and", "or", or "&".
CultureCode	String	The culture codes contained in the input data.
CultureCodeUsedToParse	String	Identifies the culture-specific grammar that was used to parse the data. Null (empty) Global culture (default). de German. es Spanish. ja Japanese. Note: If you added your own domain using the Open Parser Domain Editor, the cultures and culture codes for that domain will appear in this field as well.
FirstName	String	The first name of a person.
GeneralSuffix	String	A person's general/professional suffix. For example, MD or PhD.
IsParsed	String	Indicates whether an output record was parsed. Values are true or false.
IsPersonal	String	Indicates whether the name is an individual rather than a firm. Values are true or false.
IsReverseOrder	String	Indicates whether the input name is in reverse order. Values are true or false.
LastName	String	The last name of a person. Includes the paternal last name.

Field Name	Format	Description
LeadingData	String	Non-name information that appears before a name.
MaturitySuffix	String	A person's maturity/generational suffix. For example, Jr. or Sr.
MiddleName	String	The middle name of a person.
Name.	String	The personal or firm name that was provided in the input.
NameScore	String	Indicates the average score of known and unknown tokens for each name. The value of NameScore will be between 0 and 100, as defined in the parsing grammar. 0 is returned when no matches are returned.
SecondaryLastName	String	In Spanish parsing grammar, the surname of a person's mother.
TitleOfRespect	String	Information that appears before a name, such as "Mr.", "Mrs.", or "Dr."
TrailingData	String	Non-name information that appears after a name.
Fields Related to Conjoined Names		
Conjunction2	String	Indicates that a second, conjoined name contains a conjunction such as "and", "or", or "&".
Conjunction3	String	Indicates that a third, conjoined name contains a conjunction such as "and", "or", or "&".
FirmName2	String	The name of a second, conjoined company. For example, Baltimore Gas & Electric dba Constellation Energy.
FirmSuffix2	String	The suffix of a second, conjoined company.

Field Name	Format	Description
FirstName2	String	The first name of a second, conjoined name.
FirstName3	String	The first name of a third, conjoined name.
GeneralSuffix2	String	The general/professional suffix for a second, conjoined name. For example, MD or PhD.
GeneralSuffix3	String	The general/professional suffix for a third, conjoined name. For example, MD or PhD.
IsConjoined	String	Indicates that the input name is conjoined. An example of a conjoined name is "John and Jane Smith."
LastName2	String	The last name of a second, conjoined name.
LastName3	String	The last name of a third, conjoined name.
MaturitySuffix2	String	The maturity/generational suffix for a second, conjoined name. For example, Jr. or Sr.
MaturitySuffix3	String	The maturity/generational suffix for a third, conjoined name. For example, Jr. or Sr.
MiddleName2	String	The middle name of a second, conjoined name.
MiddleName3	String	The middle name of a third, conjoined name.
TitleOfRespect2	String	Information that appears before a second, conjoined name, such as "Mr.", "Mrs.", or "Dr."

Field Name	Format	Description
TitleOfRespect3	String	Information that appears before a third, conjoined name, such as "Mr.", "Mrs.", or "Dr."

Open Name Parser Summary Report

The Open Name Parser Summary Report lists summary statistics about the job, such as the total number of input records and the total number of records that contained no name data, as well as several parsing statistics. For instructions on how to use reports, see the *Spectrum Technology Platform Dataflow Designer's Guide*.

General Results

- **Total number of input records**—The number of records in the input file.
- **Total number of records that contained no name data**—The number of records in the input file that did not contain name data to be parsed.
- **Total number of names parsed out**—The number of names in the input file that were parsed.
- **Total Records**—The total number of records processed.
- **Lowest name parsing score**—The lowest parsing score given to any name in the input file.
- **Highest name parsing score**—The highest parsing score given to any name in the input file.
- **Average name parsing score**—The average parsing score given among all parsed names in the input file.

Personal Name Parsing Results

- **Number of personal name records written**—The number of personal names in the input file.
- **Number of names parsed from conjoined names**—The number of parsed names from records that contained conjoined names. For example, if your input file had five records with two conjoined names and seven records with three conjoined names, this value for this field would be 31, as expressed in this equation: $(5 \times 2) + (7 \times 3)$.
- **Records with 2 conjoined names**—The number of input records containing two conjoined names.
- **Records with 3 conjoined names**—The number of input records containing three conjoined names.
- **Number of names with title of respect present**—The number of parsed names containing a title of respect.
- **Number of names with maturity suffix present**—The number of parsed names containing a maturity suffix.
- **Number of names with general suffix present**—The number of parsed names containing a general suffix.

- **Number of names that contained account descriptions**—The number of parsed names containing an account description.
- **Total Reverse Order Names**—The number of parsed names in the reverse order, resulting in the output field isReversed as "True".

Business Name Parsing Results

- **Number of business name records written**—The number of business names in the input file.
- **Number of names with firm suffix present**—The number of parsed names containing a firm suffix.
- **Number of names that contained account descriptions**—The number of input records containing an account description.
- **Total DBA Records**—The number of input records containing Doing Business As (DBA) conjunctions, resulting in both output fields isPersonal and isFirm as "True".

10 - Accessing Stages through Flow Designer

In this section

Advanced Matching.....380



Advanced Matching

Filter

The **Filter** stage retains or removes records from a group of records based on the rules you specify.

To use the **Filter** stage, perform these steps:

1. On the **Flow Designer** Home page, click **New**.
2. On the **New Flow** page, click **Job**, **Service**, or **Subflow**, as required and then click the corresponding blank canvas.
3. Click **OK**.
4. In the dialog box that appears, enter a name for the **Flow**, **Job**, **Service**, or **Subflow** you are creating.
5. Click **OK**.
6. On the **Stages** panel, expand the **Advanced Matching** panel and drag the **Filter** stage to the canvas.
7. Connect the required **Sources** and **Sinks** to the **Filter** stage.
8. Click **Filter**. The **Filter Options: Filter** page is displayed. Use this page to:
 - a. Define grouping and sorting options for the records
 - b. Configure filtering rules

Configuring the Filter Options

Configuring the grouping and sorting options

On the **Filter Options: Filter** page, specify these details:

Field	Description
Group by	<p>Specifies the field to use to create groups of records to filter. The Filter stage will retain one or more records from each group, depending on how you configure the stage. In cases where you have used a matching stage earlier in the dataflow, such as Interflow Match, Intraflow Match, or Transactional Match, you should select the CollectionNumber field to use the collections created by the matching stage as the groups. However, if you want to group records by some other field, choose the field here. For example, if you want to filter out all but one record from records that have the same value in the AccountNumber field, you would select AccountNumber.</p>
Sort	<p>If you specify a field in the Group by field, check this box to sort the records by the value in the field you chose. This option is enabled by default.</p>

Field	Description
Advanced	<p>Click this button to specify sort performance options. The sort performance options specified in Management Console, which are the default performance options for your system and are in effect. If you want to override your system's default performance options, toggle the Override sort performance options to yes and specify the values you want in these fields:</p> <p>In memory record limit Specifies the maximum number of data rows a sorter will hold in memory before it starts paging to disk. By default, a sort of 10,000 records or less will be done in memory and a sort of more than 10,000 records will be performed as a disk sort. The maximum limit is 100,000 records. An in-memory sort is much faster than a disk sort, so this value should be set high enough so that most of the sorts are in-memory sorts and only large sets are written to disk.</p> <p style="text-align: center;">Note: Be careful in environments where there are jobs running concurrently because increasing the In memory record limit setting increases the likelihood of running out of memory.</p> <p>Maximum number of temporary files Specifies the maximum number of temporary files that may be used by a sort process. Using a larger number of temporary files can result in better performance. However, the optimal number is highly dependent on the configuration of the server running Spectrum Technology Platform. You should experiment with different settings, observing the effect on performance of using more or fewer temporary files. To calculate the approximate number of temporary files that may be needed, use this equation:</p> <div style="border: 1px solid #ccc; padding: 5px; margin: 10px 0;"> $\frac{(NumberOfRecords \times 2)}{InMemoryRecordLimit} = NumberOfTempFiles$ </div> <p>Note that the maximum number of temporary files cannot be more than 1,000.</p> <p>Compression Specifies that temporary files are compressed when they are written to disk.</p> <p style="text-align: center;">Note: The optimal sort performance settings depends on your server's hardware configuration. You can use this equation as a general guideline to produce good sort performance: $(InMemoryRecordLimit \times MaxNumberOfTempFiles \div 2) \geq TotalNumberOfRecords$</p>
Limit number of returned duplicate records	<p>Specifies the maximum number of records that are returned from each group. If you set this option to 1, you can define filter rules to determine which record in each group should be returned. If no rules are defined, the first record in each collection is returned and the rest are discarded. In this mode, the filter rules define which record will be retained.</p> <p>For example, if you define a rule where the record with the highest match score in a group is retained, and you set this option to 1, then the record with the highest match score in each group will survive and the other records in the group will be discarded.</p> <p>If you set this option to a value higher than one, you cannot specify filter rules.</p> <p style="text-align: center;">Note: In the event no records in the collection meet the defined rule criteria, then no records from the group are returned.</p>

Field	Description
Remove duplicates from collection	<p>Specifies to use filter rules to determine which records are removed from the collection. The remaining records in the collection are retained. When this option is selected, you must define a rule.</p> <p>Note: If a group contains only one record, the filter rules are ignored and the record is retained.</p>

Configuring the filter rules

Filter rules determine which records in a group are to be retained or removed. If you select the option **Limit number of returned duplicate records** then the rules determine which records survive the filter. If you select the option **Remove duplicates from collection** then the rules determine which records are removed from the dataflow.

1. To configure filter rules, click the **Add Rule**  icon in the **Rule Configuration** section, and enter details as described in the table below.

Field	Description
Field name	Specifies the name of the dataflow field whose value you want to evaluate to determine whether to filter the record.
Field Type	<p>Specifies the type of data in the field. It can be one of these:</p> <ul style="list-style-type: none"> • Non-Numeric: Choose this option if the field contains non-numeric data (for example, string data). • Numeric: Choose this option if the field contains numeric data (for example, double, and float).

Field	Description
Operator	<p>Specifies the type of comparison you want to use to evaluate the field. Use one of these options:</p> <ul style="list-style-type: none"> • Contains: Determines if the field contains the value specified. For example, "sailboat" contains the value "boat". • Equal: Determines if the field contains the exact value specified. • Greater Than: Determines if the field value is greater than the value specified. • Greater Than Or Equal To: Determines if the field value is greater than or equal to the value specified. • Highest: Compares the field's value for all the records group and determines which record has the highest value in the field. For example, if the fields in the group contain values of 10, 20, 30, and 100, the record with the field value 100 would be selected. If multiple records are tied for the longest value, one record is selected. • Is Empty: Determines if the field contains no value. • Is Not Empty: Determines if the field contains any value. • Less Than: Determines if the field value is less than the value specified. • Less Than Or Equal To: Determines if the field value is less than or equal to the value specified. • Longest: Compares the field's value for all the records group and determines which record has the longest (in bytes) value in the field. For example, if the group contains the values "Mike" and "Michael", the record with the value "Michael" is selected. If multiple records are tied for the longest value, one record is selected. • Shortest: Compares the field's value for all the records group and determines which record has the shortest (in bytes) value in the field. For example, if the group contains the values "Mike" and "Michael", the record with the value "Mike" is selected. If multiple records are tied for the shortest value, one record is selected. • Lowest: Compares the field's value for all the records group and determines which record has the lowest value in the field. For example, if the fields in the group contain values of 10, 20, 30, and 100, the record with the field value 10 would be selected. If multiple records are tied for the lowest value, one record is selected. • Most Common: Determines if the field value contains the value that occurs most frequently in this field among the records in the group. If two or more values are most common, no action is taken. • Not Equal: Determines if the field value is not the same as the value specified.

Field	Description
Value type	<p>Specifies the type of value you want to compare to the field's value. Use one of these values:</p> <p>Note: This option is not available if you select the operators <i>Highest, is Empty, is Not empty, Lowest, Longest, Most Common and Shortest</i>.</p> <ul style="list-style-type: none"> • Field: Choose this option if you want to compare another dataflow field's value to the field. • String: Choose this option if you want to compare the field to a specific value.
Value	<p>Specifies the value to compare to the field's value. If you selected Field in the Field type field, select a dataflow field. If you selected String in the Value type field, type the value you want to use in the comparison.</p> <p>Note: This option is not available if you select the operators <i>Highest, is Empty, is Not empty, Lowest, Longest, Most Common and Shortest</i>.</p>

- To define another rule, click the **Add Rule**  icon and specify details as described in the table above.

Note: If you specify multiple rules, you will have to select a logical operator between each rule. Choose **And** if you want the new rule and the previous rule to both pass in order for the condition to be met. Select **Or** if you want either the previous rule or the new rule to pass in order for the condition to be met.

Example of a Filter Rule

This rule retains the record in each group with the highest value in the MatchScore field. Note that **Value** and **Value Type** options do not apply when the Operator is highest or lowest.

Field Name = MatchScore
 Field Type = Numeric
 Operator = Highest

This rule retains the record where the value in the AccountNumber is "12345".

Field Name = AccountNumber
 Field Type = Numeric
 Operator = Equals
 Value Type = String
 Value = 12345

Match Key Generator

This stage creates a match key for each record. A match key is a non-unique key shared by like records that identify records as potential duplicates. The match key facilitates the matching process by only comparing records that contain the same match key.

This stage allows you to:

- **Create a match key:** Select **Create Match Key** at the left corner of the page to create a new match key. It also allows you to publish a match key to the repository for re-use at a later point in time. For more information, see [Creating a New Match Key](#) on page 387.
- **Preview a match key:** You can preview the match key with the **Input Data** and the **Result** at the bottom of the **Options** page to avoid any potential errors.
- **Configure access control for match key:** It allows you to regulate who can view your match key and define the permissions for them. For more information, see [Configuring Access Control for Match Key](#) on page 391.

How Match Key Generator functions

The match key is created using rules you define and is comprised of multiple rules. Each rule has a selected algorithm that is applied on the chosen input field. The result of each rule is then linked together to create a single match key.

Example of a Match Key Generator Rule

Consider the input as below:

- First Name - Fred
- Last Name - Mertz
- Gender Code - M

With the above input, if you are required to define a match key rule that generates a match key by combining data from the record shown below, the match key is MertzFredM.

Input Field	Start Position	Length
Last Name	1	5
First Name	1	5
Gender Code	1	1

Creating a New Match Key

The match key generator allows you to define rules for your match key with various options depending on your requirements. You can define any number of rules for a single match key.

To create a new match key, perform these steps:

1. On the **Match Key Generator Options** page, select **Create Match Key** from the **Match Keys** panel located on the left corner of the page.
2. (Optional) If you want to create a new match key by using an existing match key stored in the repository with predefined rules, select **Copy From Repository**. For more information, see [Creating a Match Key from Repository](#) on page 390.
3. In the **Match key name** field, enter a unique name for the match key you want to create.
4. To define the rule for your match key, select **Add**.
You see the **Rule Options** panel on the right corner of the page.
5. In the **Rule Options** panel, configure these options:

Option	Description	
Algorithm	Consonant	It returns specified fields with consonants removed.
	Double Metaphone	It returns a code based on a phonetic representation of their characters. Double Metaphone is an improved version of the metaphone algorithm and attempts to account for the many irregularities found in different languages.
	Koeln	Indexes names by sound as they are pronounced in German. It allows names with the same pronunciation to be encoded to Koeln the same representation so that they can be matched, despite minor differences in spelling. The result is always a sequence of numbers; special characters and white spaces are ignored. This option was developed to respond to the limitations of Soundex.
	MD5	It is a message-digest algorithm that produces a 128-bit hash value, which is used to check data integrity.
	Metaphone	It returns a Metaphone coded key of selected fields. Metaphone is an algorithm for coding words using their English pronunciation.
	Metaphone (Spanish)	It returns a Metaphone coded key of selected fields for the Spanish language. This Metaphone algorithm codes words using their Spanish pronunciation.
	Metaphone3	It improves upon the Metaphone and Double Metaphone algorithms with a more exact consonant and internal vowel Metaphone3 settings that allow you to produce words or names more or less closely matched to search terms on a phonetic basis. Metaphone3 increases the accuracy of phonetic encoding to 98%. This option was developed to respond to the limitations of Soundex.
	Nysiis	It is a phonetic code algorithm that matches an approximate pronunciation to an exact spelling and indexes words that are Nysiis pronounced similarly. It is a part of the New York State Identification and Intelligence System. For example, you are looking for someone's information in a database of people. You believe that the person's name sounds like "John Smith," but it is, in fact, spelled "Jon Smyth." If you conducted a search looking for an exact match for "John Smith," no results would be returned. However, if you index the database using the NYSIIS algorithm and search using the NYSIIS algorithm again, the correct match will be returned because both "John Smith" and "Jon Smyth" are indexed as "JAN SNATH" by the algorithm.
	Phonix	It preprocesses name strings by applying more than 100 transformation rules to single characters or to sequences of several characters. Nineteen of those rules are applied only if the characters are at the beginning of the string, while 12 of the rules are applied if they are in the middle of the string, and 28 of the rules are applied if they are at the end of the string. The transformed name string is encoded into a code that is comprised by a starting letter followed by three digits (removing zeros and duplicate numbers). This option was developed to respond to the limitations of Soundex; it is more complex and, therefore, slower than Soundex.
	Sonnex	This algorithm determines the similarity between 2 French-language strings based on the phonetic representation of their characters. It returns a Sonnex coded key of the selected fields.
	Soundex	It returns a Soundex code of selected fields. Soundex produces a fixed-length code based on the English pronunciation of a word.
Substring	It returns a specified portion of the selected field.	

Option	Description	
Field name	Select the field to which you want to apply the selected algorithm. For example, if you select LastName , and you choose the SubString algorithm, the SubString algorithm is applied to the data in the LastName field to produce a match key.	
Start Position	Select the start position within the specified field. Note: You can select a start position only if you choose the SubString algorithm.	
Length	Select the length of characters you want to include from the starting position. Note: You can select a length only if you choose the SubString algorithm.	
Pre-processing options	Remove noise character	It removes all non-numeric and non-alpha characters such as hyphens, white space, and other special characters from an input field.
	Sort input	It sorts all characters or terms in an input field in alphabetical order. <ul style="list-style-type: none"> • Characters: It sorts the values of the characters from an input field prior to creating a unique ID. • Terms: It sorts each term value from an input field prior to creating a unique ID. <p>Note: You can select the sort input only if you choose the Consonant or SubString algorithm.</p>

6. Select **OK**.
The rule you define for the match key appears below **Rules** on the **Create Match Key** workspace.
7. (Optional) If you want to deactivate the rule for the match key, clear the **Active** checkbox.
8. (Optional) Select the **Modify**  button to modify the rule.
9. (Optional) Select the **Delete**  button to delete the rule.
10. (Optional) Follow steps 4 on page 387 through 5 on page 387 to define any number of rules for your match key.
11. On the right corner, select **Create**.
Your match key is created and appears on the **Match Keys** panel located on the left corner. You can see the preview of the match key at the bottom of the page.

Note: If you want to update the match key, you can select and edit the match key from the **Match Keys** panel.

12. (Optional) If you want to store your match key in the repository for re-use:
 - a. Select **Publish**.

- b. On the **Publish Match Key** panel, enter a name with which you want to store the match key in the repository. In case a match key is already stored with the same name, you can:
 - Enter a unique name in the **Match key name** field, or
 - Select the **override** checkbox to override the existing match key.
- c. Select **OK**.

Note: You can regulate the access to your match key if you do not want any other user to re-use them. For more information, see [Configuring Access Control for Match Key](#) on page 391.

13. On the top-right corner of the page, select **Apply**.
Your match key applies to the match key generator stage in the dataflow.

Creating a Match Key from Repository

The **Copy From Repository** option allows you to save significant time by re-using the published match keys from the repository. You can use the match key as-is or modify it to create a new match key.

Note: You can view or modify an existing match key from the repository only if you have the required user privileges. For more information, see [Configuring Access Control for Match Key](#) on page 391.

To create a match key from the repository, perform these steps:

1. In the **Create Match Key** workspace, select **Copy From Repository**.
2. In the **Match Key Repository** field, select an existing match key with pre-defined rules from the drop-down list.
3. Select **OK**.
The match key appears in the **Create Match Key** workspace.

Note: If you want to create a match key as-is, skip steps **4** on page 390 and **5** on page 390.

4. (Optional) If you want to enter a new name for the match key, use the **Match key name** field.
5. (Optional) If you want to add new rules to the match key, select **Add**.
6. On the right corner, select **Create**.
Your match key is created and appears on the **Match Keys** panel located on the left corner. You can see the preview of the match key at the bottom of the **Match Key Generator Options** page.
7. (Optional) If you want to store your match key in the repository, publish the match key.
8. On the top-right corner of the page, select **Apply**.
Your match key applies to the match key generator stage in the dataflow.

Configuring Access Control for Match Key

Access control to a match key minimizes the risk of unauthorized usage by a user. It allows you to regulate who can view your match key and define the permissions for your match key.

To configure access control for your match key, perform these steps:

1. Go to **Management Console > System > Security**.
2. Select the **Access Control** tab.
3. Select the **Add access control**  button.
4. Do one of these:
 - If you want to specify access controls for a role, select **Role** from the drop-down field. The access control permissions you specify will affect all users who have the role you choose.
 - If you want to specify access controls for a single user, select **User** from the drop-down field. The access control permissions you specify will only affect the user you choose.
5. Select the **Add entity override**  button. You see an **Add Entity** workspace.
6. In the **Add Entity** workspace, select **AMM.Match Keys Management** from the drop-down field and do these:
 - a. Drag the match keys you see on the left side to the **Selected Entities** space on the right side.
 - b. Select **Add**.

The match keys you add are displayed in the **AMM.Match Keys Management** panel at the bottom on the **Add Access Control** page.
7. Specify the **View**, **Modify**, and **Delete** permissions that you want to grant for each match key. You can select one of these options for each permission:

Option	Description
	The permission is inherited from the role.
	The permission is inherited from the role and cannot be overridden.
	The permission is granted, overriding the permission specified in the user or role.
	The permission is denied, overriding the permission specified in the user or role.

8. On the top-right corner of the page, select **Save**.
Your configuraton applies to the match keys.

Using Command Line Interface for Match Key

You can use the CLI to perform various actions on your match keys. For more details about the CLI commands, see the Administration Utility section of the *Spectrum Technology Platform Administration Guide*.

Output

For more information, see [Output](#) on page 299.

Write to Search Index

The **Write to Search Index** stage enables you to create a full-text index based on the data coming into the stage. Having this data in a dedicated search index results in quicker response time when you conduct searches against the index from other Spectrum Technology Platform stages.

Full-text-search indexes are preferable to relational databases when you have a great deal of free-form text data that needs to be searched or categorized or if you support a high volume of interactive, text based queries.

This stage allows you to:

- **Create an Index:** In the **Index** field, select the **Create Index**  button to create an index. For more information, see [Creating an Index](#) on page 393.
- **View the schema of the fields:** The **Fields** section allows you to view the schema of the fields for the selected index. You need to configure the source fields for all the fields listed as necessary.
- **Apply Write to Search Index:** Click **Apply** at the top-right corner of the page to apply the configuration changes to the write to search index stage in the dataflow. For more information, see [Applying Configuration Changes to the DataFlow](#) on page 396.

How Write to Search Index functions

Write to Search Index uses analyzers to break input text into small indexing elements called tokens. It then extracts search index terms from those tokens. The type of analyzer used—how input text is broken into tokens—determines how you will then be able to search for that text.

For example, the Keyword analyzer always does an exact match and tokenizes the whole string as a single token, while the Standard analyzer breaks the string to create tokens. The Stop Word analyzer is more sophisticated and removes articles, such as "a" or "the" from the string before tokenizing, thus shrinking the index size. Search indexes support the near real-time feature, allowing indexes to be updated almost immediately, without the need to close and rebuild the stage using the search index.

Creating an Index

The write to search index stage allows you to select the fields you want to be stored and indexed from your input source.

To create an index, perform these steps:

1. On the **Write to Search Index Options** page, select the **Create Index**  button in the **Index** field.
You see a **Create Index** panel on the right corner of the page.
2. In the **Index name** field, enter a unique name for the index.

Note: You must enter the name of the index as alphanumeric only, for example, xyz123 or xyz, and it must not contain any special characters. Also, it is mandatory to enter a name. If you do not enter a name, you cannot create an index.

3. To select the fields you want to store and index from your input source, go to the **Fields** table and do as follows:

Note: It is mandatory to index and store at least one field. You can store and index the same field, or you can index one field and choose to store a different field to create an index.

- a) (Optional) Use **Filter**  to selectively add the fields you want to index from the input source.
- b) (Optional) In the **Field name**, if you want to change the name of the field, enter the new name.

Note: You must enter the field name as alphanumeric only. Make a note that you cannot change the **Stage field** name and is disabled for your use.

- c) Select the field **Type** - **double**, **float**, **integer**, **long**, or **string**.
- d) To store the field, select the **Store** checkbox.
- e) To add the field to the index for search query, select the **Index** checkbox.
You see that the **Analyzer** is enabled.
- f) In the **Analyzer** field, select one of these analyzers you want to use to index the field:

Analyzer	Description
Danish	It supports Danish-language indexes and type-ahead services. It also supports many stop words and removes articles such as "at," "and," and "a" to shrink the index size and increase performance.
Dutch	It supports Dutch-language indexes and type-ahead services. It also supports many stop words and removes articles such as "the," "and," and "a" to shrink the index size and increase performance.

Analyzer	Description
Finnish	It supports Finnish-language indexes and type-ahead services. It also supports many stop words and removes articles such as "is," "and," and "of" to shrink the index size and increase performance.
French	It supports French-language indexes and type-ahead services. It also supports many stop words and removes articles such as "the," "and," and "a" to shrink the index size and increase performance.
German	It supports German-language indexes and type-ahead services. It also supports many stop words and removes articles such as "the," "and," and "a" to shrink the index size and increase performance.
Hindi	It supports Hindi-language indexes and type-ahead services. It also supports many stop words and removes articles such as "by," "and," and "a" to shrink the index size and increase performance.
Hungarian	It supports Hungarian-language indexes and type-ahead services. It also supports many stop words and removes articles such as "the," "and," and "a" to shrink the index size and increase performance.
Italian	It supports Italian-language indexes and type-ahead services. It also supports many stop words and removes articles such as "the," "and," and "a" to shrink the index size and increase performance.
Keyword	It creates a single token from a stream of data and keeps it as is. For example, the string 'XYZ Software' would be returned as just one token 'XYZ Software.' For a comparison of Standard and Keyword analyzers, see Standard and Keyword Analyzer on page 309.
Norwegian	It supports Norwegian-language indexes and type-ahead services. It also supports many stop words and removes articles such as "the," "and," and "a" to shrink the index size and increase performance.
Portugese	It supports Portuguese-language indexes and type-ahead services. It also supports many stop words and removes articles such as "the," "and," and "a" to shrink the index size and increase performance.
Russian	It supports Russian-language indexes and type-ahead services. It also supports many stop words and removes articles such as "and," "I," and "you" to shrink the index size and increase performance.
Spanish	It supports Spanish-language indexes and type-ahead services. It also supports many stop words and removes articles such as "the," "and," and "a" to shrink the index size and increase performance.

Analyzer	Description
Standard	It provides a grammar based tokenizer that contains a superset of the Whitespace and Stop Word analyzers. It understands English punctuation for breaking down words, knows words to ignore (via the Stop Word Analyzer), and performs case-insensitive searching technically by conducting lowercase comparisons. For example, the string 'XYZ Software' would be returned as two tokens: 'XYZ,' and 'Software.' For a comparison of Standard and Keyword analyzers, see Standard and Keyword Analyzer on page 309.
StopWord	It removes articles such as "the," "and," and "a" to shrink the index size and increase performance.
Swedish	It supports Swedish-language indexes and type-ahead services. It also supports many stop words and removes articles such as "the," "and," and "a" to shrink the index size and increase performance.
Whitespace	It separates tokens with whitespace. Somewhat of a subset of the Standard Analyzer in that it understands word breaks in English text based on spaces and line breaks.

Note: The **Analyzer** is enabled only if you choose to index the field. Also, it is mandatory to select the analyzer for each index field.

4. Select a key field from the **Key field** drop-down list on the basis of which you want to **Write mode: Delete or Update or append** the records. See the example below.

In case of **Create or overwrite** mode, the **Key field** needs to be unique for Elastic search indexes (used in a distributed environment). The table below explains the indexing behavior if the **Key field** is non-unique for Elastic search index.

Write mode	Key field	Elastic search index
Create or Overwrite	Duplicate records with same Key field	All duplicate records with the same Key field are overwritten .
Update or Append	Duplicate records with same Key field	Duplicates are overwritten.

Note: You see a key field in the drop-down list only if you choose to index the field from the **Fields** table.

5. Select **Create**.
The index you create appears on the **Index** drop-down list on the **Write to Search Index Options** page.

Applying Configuration Changes to the DataFlow

To apply the configuration changes to the dataflow, perform these steps on the **Write to Search Index Options** page:

1. In the **Index** field, select the index for which you want to use the schema for the current stage configuration from the drop-down list.
2. In the **Write Mode** field, select one of these write mode from the drop-down list that gives you the options related to how the new data should affect the existing data:

Write Mode	Description
Create or overwrite	The new data overwrites the existing data, and the existing data remains no longer in the index.
Update or append	The new data overwrites the existing data, and any new data that did not previously exist is added to the index. It also allows you to update the blank fields by selecting the Update blank fields option. If the input is blank and you do not want the index field to update with blank data, you can deselect the Update blank fields option.
Append	It adds new data to the existing data, and the existing data remains intact.
Delete	It deletes the data for the selected field from the search index.

3. (Optional) Select **Batch commit** to enable the batch option. Use the **Batch size** option to define the number of records you want to commit in one go. The value available by default is 5000.
4. On the top-right corner of the page, select **Apply**.
Your configuration changes applies to the write to search index stage in the dataflow.

Output

For more information, see [Output](#) on page 308.

Intraflow Match

The **Intraflow Match** stage locates matches between similar data records within a single input stream. You can create hierarchical rules based on any fields that have been defined or created in other stages of the dataflow.

This stage allows you to:

- Create and publish a new match rule to the repository for re-use at a later point in time, which allows you to save significant time. For more information, see [Creating a match rule](#) on page 400.
- Use a predefined match rule as-is from **Template Rules** in **Rule Configuration** panel. It gives you the flexibility to choose a sophisticated set of pre-defined rules that saves time than defining a rule from scratch.
- Configure the grouping and sorting options for the match rule in the **Settings** panel. For more information, see [Configuring the grouping and sorting options](#) on page 397.

How Intraflow Match functions

Using **Group by** set by you, the matcher identifies groups of records that might potentially be duplicates of one another. The matcher then proceeds through each record in the group.

- If the record matches an existing suspect, the record is considered a duplicate of that suspect, assigned a score, collection number, and match record type (Duplicate), and eliminated from the match.
- If the record matches no existing suspect within the match group, the record becomes a new suspect, in that it is added to the current match group so that it can be matched against by subsequent records.

When the matcher exhausts all records in the current match group, it eliminates all suspects from the match, labeling the match record type as unique and assigning a collection number of 0. Those suspects with a least one duplicate will retain a match record type of suspect and are assigned the same collection number as its matched duplicate record. Finally, when all records within a match group have been written to the output, a new match group is compared.

Note: The default matching method will only compare records that are within the same match group.

The type of matching (Intraflow or Interflow) determines how express key match results translate to candidate match scores.

- In Intraflow matching, the score a candidate gains as a result of an express key match depends on whether the record to which that candidate matched was a match of some other suspect—express key duplicates of a suspect will always have match scores of 100. In contrast, express key duplicates of another candidate (which was a duplicate of a suspect) will inherit the match score (not necessarily 100) of that candidate.
- In Interflow matching, a successful express key match always confers a 100 match score onto the candidate.

Configuring Options for Intraflow Match

Configuring the grouping and sorting options

On the **Intraflow Match Options: Intraflow Match** page, configure these options in the **Settings** panel:

Option	Description
Group by	It selects a field to use for grouping records in the match queue. Intraflow Match only attempts to match records against other records in the same match queue.
Sort	It performs a pre-match sort of your input based on the selection in the Group by field.
Advanced	Select this button to specify additional sort options. If you want to override your system's default performance options, toggle the Override sort performance options to ON and specify the values you want in these fields: <div data-bbox="492 674 1414 877"> <p>In memory record limit Specifies the maximum number of data rows a sorter will hold in memory before it starts paging to disk. By default, a sort of 10,000 records or less will be done in memory, and a sort of more than 10,000 records will be performed as a disk sort. The maximum limit is 100,000 records. An in-memory sort is much faster than a disk sort, so this value should be set high enough so that most of the sorts are in-memory sorts, and only large sets are written to disk.</p> <p>Note: Be careful in environments where jobs are running concurrently because increasing the In memory record limit setting increases the likelihood of running out of memory.</p> </div> <div data-bbox="492 1024 1414 1260"> <p>Maximum number of temporary files Specifies the maximum number of temporary files that may be used by a sort process. Using a larger number of temporary files can result in better performance. However, the optimal number is highly dependent on the configuration of the server running Spectrum Technology Platform. You should experiment with different settings, observing the effect on the performance of using more or fewer temporary files. To calculate the approximate number of temporary files that may be needed, use this equation:</p> <div data-bbox="670 1283 1414 1371" style="background-color: #f0f0f0; padding: 5px;"> $\frac{(NumberOfRecords \times 2)}{InMemoryRecordLimit} = NumberOfTempFiles$ </div> <p>Note: The maximum number of temporary files cannot be more than 1,000.</p> </div> <div data-bbox="492 1486 1414 1686"> <p>Compression Specifies that temporary files are compressed when they are written to disk.</p> <p>Note: The optimal sort of performance settings depends on your server's hardware configuration. You can use this equation as a general guideline to produce good sort performance: $(InMemoryRecordLimit \times MaxNumberOfTempFiles \div 2) >= TotalNumberOfRecords$</p> </div>
Express match on	It performs an initial comparison of express key values to determine whether two records are considered a match.

Option	Description
--------	-------------

Initial collection number

It specifies the starting number to assign to the collection number field for duplicate records. The collection number identifies each duplicate record in a match queue. Unique records are assigned a collection number of 0. Each duplicate record is assigned a collection number starting with the value specified in the **Initial collection number** field.

Assign collection number 0 to unique records

It assigns zeroes as collection numbers to unique records. Un-check this option to generate collection numbers other than zero for unique records. The unique record collection numbers will be in sequence with any other collection numbers. See the examples below.

- If your matching dataflow finds four records and the first two records are unique, the collection numbers would be assigned, as shown in the table below:

Collection Number	Record Type
1	Unique
2	Unique
3	Duplicate or Suspect
3	Duplicate or Suspect

- If your matching dataflow finds four records and the last two are unique, the collection numbers would be assigned, as shown in the table below:

Collection Number	Record Type
1	Duplicate or Suspect
1	Duplicate or Suspect
2	Unique
3	Unique

Option	Description
Generate data for analysis	It generates the match results. For more information, see Analyzing Match Results on page 111.
Return match rule name	It includes the selected match rule name in the stage output.
Enable sliding window	<p>It enables the sliding window for this matching method and allows you to enter a value in the Window size field. By default, the value is 50.</p> <p>The sliding window algorithm sequentially fills a predetermined buffer size called a window with the corresponding amount of data rows. As each row is added to the window, it's compared to each item already contained in the window. If a match with an item is determined, then the same group ID is given to:</p> <ul style="list-style-type: none"> • the driver record (the new item to add to the window), and • the candidates (items already in the window). <p>This comparison is continued until the driver record has been compared to all items contained within the window.</p> <p>As new drivers are added, the window will eventually reach its predetermined capacity. At this point, the window will slide, hence the term sliding window. Sliding means that the window buffer will remove and write the oldest item in the window as it adds the newest driver record to the window.</p>

Configuring the match rules

On the **Rule Configuration** panel, you can:

- Configure one of the predefined rules from **Template Rules** located on the left corner of the page which you can use as-is, or
- Configure a new match rule and publish it to the repository for re-use. You can configure a new match rule without using one of the predefined match rules. For more information, see [Creating a match rule](#) on page 400.

Creating a match rule

To create a new match rule and publish it to the repository for re-use, perform these steps:

1. On the **Rules** panel, click **Create new rule**.

Note: By default, Custom is added as a prefix to the **Rule Name** followed by the name you define in Step 2 on page 400, for example, Custom - xyz.

2. Specify the dataflow fields, parent or child; you want to use in the match rule and match rule hierarchy.

- Click  button and enter a name for the parent below **Match when not true**.

Note: The name you enter for your first parent in the hierarchy is used as the match rule name, for example, Custom - xyz.

- Click  button and select a field to add to the parent from the drop-down list below **Match when not true**.

Note: All children under a parent must use the same logical operator. If you want to use different logical operators between fields, you must first create intermediate parents.

3. Define these parent options as listed in the table below, which are displayed on the parent node:

Option	Description
Match when not true	It changes the logical operator for the parent from and to and not . If you select this option, records will only match if they do not match the logic defined in this parent. Note: If you select the Match when not true option, it negates the Matching Method options. For more information, see Negative Match Conditions on page 77.
Matching Method	Select one of these from the drop-down list to determine if a parent is a match or non-match: <ul style="list-style-type: none"> • All true: A parent is considered a match if all children are determined to match. This method creates an "and" connector between children. • Any true: A parent is considered a match if at least one child is determined to match. This method creates an "or" connector between children. • Based on threshold: A parent is considered a match if at least one child is determined to match. This method creates an "or" connector between children. <p>If you select this option, the Threshold field enables you to specify a threshold value. The Scoring Method determines which logical connector to use. The thresholds at the parent cannot be higher than the threshold of the children. For more information, see the matching method-to-scoring method matrix below this table.</p>
Missing Data	Select one of these from the drop-down list to specify how to score blank data in a field: <ul style="list-style-type: none"> • Ignore blanks: Ignores the field if it contains blank data. • Count as 0: Scores the field as 0 if it contains blank data. • Count as 100: Scores the field as 100 if it contains blank data. • Compare blanks: Scores the suspect and candidate fields as 100 if they both contain blank data; otherwise, scores the suspect and candidate fields as 0.

Option	Description
Scoring Method	<p>Select one of these from the drop-down list to determine the matching score:</p> <ul style="list-style-type: none"> • Weighted Average: Uses the weight of each child to determine the average match score. • Average: Uses the average score of each child to determine the score of a parent. • Maximum: Uses the highest child score to determine the score of a parent. • Minimum: Uses the lowest child score to determine the score of a parent. • Vector Summation: Uses the vector summation of each child score to determine the score of the parent. The formula for calculation is: $\sqrt{a^2+b^2+c^2} / \sqrt{n}$, where a, b, and c are the scores of three children, and n is the number of children. <p>For more information, see the matching method-to-scoring method matrix below this table.</p>
Evaluate	For more information, see Evaluating a match rule on page 410.
Copy settings to	<p>It allows you to copy the same settings for any number of elements.</p> <ol style="list-style-type: none"> Use the drop-down list to select or de-select the elements. Click Apply adjacent to the Copy Settings to field to copy and apply the same settings for the selected elements. <p>Note: You can copy the parent settings to a parent element and child settings to a child element only.</p>

Matching Method-to-Scoring Method Matrix

The table below shows the logical relationship between **Matching Method** and **Scoring Method** and how each combination changes the logic used during match processing.

Scoring Method	Matching Method			Comments
	Any true	All true	Based on threshold	
Weighted Average	NA	and	and	Only available when All true or Based on threshold are selected as the Matching Method .
Average	NA	and	and	
Vector Summation	NA	and	and	

Scoring Method	Matching Method			Comments
	Any true	All true	Based on threshold	
Maximum	or	NA	or	Only available when All true or Based on threshold are selected as the Matching Method .
Minimum	or	NA	or	

4. Define these child options as listed in the table below, which are displayed on the child node:

Option	Description
Match when not true	<p>It changes the logical operator from and to not. If you select this option, the match rule will only evaluate to true if the records do not match the logic defined in this child.</p> <p>For example, if you want to identify individuals who are associated with multiple accounts, you could create a match rule that matches the name but where the account number does not match. You would use the Match when not true option for the child that matches the account number.</p>
Candidate field	Select this to map the child record field you select from the drop-down list to a field in the input file.
Cross match against	Select this to choose one or more field names from the drop-down list to match different fields to one another between two records.
Threshold	Enter the threshold that must be met at the individual field level for that field to be determined a match.
Missing Data	<p>Select one of these from the drop-down list to specify how to score blank data in a field:</p> <ul style="list-style-type: none"> • Ignore blanks: Ignores the field if it contains blank data. • Count as 0: Scores the field as 0 if it contains blank data. • Count as 100: Scores the field as 100 if it contains blank data. • Compare blanks: Scores the suspect and candidate fields as 100 if they both contain blank data; otherwise, scores the suspect and candidate fields as 0.

Option	Description
Scoring Method	<p>Select one of these from the drop-down list to determine the matching score:</p> <ul style="list-style-type: none"> • Weighted Average: Uses the weight of each algorithm to determine the average match score. • Average: Uses the average score of each algorithm to determine the match score. • Maximum: Uses the highest algorithm score to determine the match score. • Minimum: Uses the lowest algorithm score to determine the match score. • Vector Summation: Uses vector summation of the score of each algorithm to determine the match score. This scoring method is useful if you want a higher vector summation match score in one or more algorithms to get proportionately represented in the final match score. The formula for calculating the final score is: $\sqrt{a^2+b^2+c^2} / \sqrt{n}$, where a, b, and c are the scores of three different algorithms, and n is the number of algorithms used.
Evaluate	For more information, see Evaluating a match rule on page 410.
Copy Settings to	<p>It allows you to copy the same settings for any number of elements.</p> <ol style="list-style-type: none"> Use the drop-down list to select or de-select the elements. Click Apply adjacent to the Copy Settings to field to copy and apply the same settings for the selected elements. <p>Note: You can copy the parent settings to a parent element and child settings to a child element only.</p>

- To configure algorithms for your child type, click **Configure Algorithms** on the child options node to add one or more of these algorithms to determine the match in the field values:

Note: Use **Search**  to selectively configure the algorithms.

String Matching Algorithms

Acronym	<p>It determines whether a business name matches its acronym by looking for acronym data; otherwise, it creates an acronym using the first character of every word.</p> <p>For example, Internal Revenue Service and its acronym IRS would be considered a match and return a match score of 100.</p>
Character Frequency	It determines the frequency of occurrence of each character in a string and compares the overall frequencies between two strings.
Exact Match	It determines if two strings are the same.
Initials	It matches the initials for parsed personal names.

Name Variant It determines whether two names are variants of each other. The algorithm returns a match score of 100 if two names are variations of each other, and a match score of 0 if two names are not variations of each other.

For example, JOHN is a variation of JAKE and returns a match score of 100. JOHN is not a variant of HENRY and returns a match score of 0.

Click **Edit** to specify the name variant options. For more information, see [Name Variant Finder](#) on page 366.

Numeric String It compares address lines by separating the numerical attributes of an address line from the characters. See the examples below.

- In the string address 1234 Main Street Apt 567, the numerical attributes of the string (1234567) are parsed and handled differently from the remaining string value (Main Street Apt). The algorithm first matches numeric data in the string with the numeric algorithm. If the numeric data match is 100, the alphabetic data is matched using Edit distance and Character Frequency. The final match score is calculated as follows:

$$(numericScore + (EditDistanceScore + CharacterFrequencyScore) / 2) / 2$$

- If you calculate the match score of these two addresses:

123 Main St Apt 567
123 Maon St Apt 567

the match score would be 95.5, calculated as follows:

Numeric Score = 100
Edit Distance = 91
Character Frequency = 91

$91 + 91 = 182$
 $182 / 2 = 91$
 $100 + 91 = 191$
 $191 / 2 = 95.5$

SubString It determines whether one string occurs within another.

Phonetic Algorithms

Daitch-Mokotoff Soundex A Phonetic algorithm that allows greater accuracy in matching of Slavic and Yiddish surnames with similar pronunciation but differences in spelling. Coded names are six digits long, and multiple possible encodings can be returned for a single name. This option was developed to respond to the limitations of Soundex in the processing of Germanic or Slavic surnames.

Double Metaphone It determines the similarity between two strings based on a phonetic representation of their characters. Double Metaphone is an improved version

of the Metaphone algorithm and attempts to account for the many irregularities found in different languages.

Koeln	Indexes names by sound as they are pronounced in German. Allows names with the same pronunciation to be encoded to the same representation so that they can be matched, despite minor differences in spelling. The result is always a sequence of numbers; special characters and white spaces are ignored. This option was developed to respond to the limitations of Soundex.
Metaphone	It determines the similarity between two English-language strings based on a phonetic representation of their characters. This option was developed to respond to the limitations of Soundex.
Metaphone (Spanish)	It determines the similarity between two strings based on a phonetic representation of their characters. This option was developed to respond to the limitations of Soundex.
Metaphone3	It improves upon the Metaphone and Double Metaphone algorithms with a more exact consonant and internal vowel settings that allow you to produce words or names more or less closely matched to search terms on a phonetic basis. Metaphone3 increases the accuracy of phonetic encoding to 98%. This option was developed to respond to the limitations of Soundex.
Nysiis	<p>It is a phonetic code algorithm that matches an approximate pronunciation to an exact spelling and indexes words that are pronounced similarly—part of the New York State Identification and Intelligence System.</p> <p>For example, consider that you are looking for someone's information in a database of people. You believe that the person's name sounds like "John Smith," but it is spelled "Jon Smath". If you conducted a search looking for an exact match for "John Smith," no results would be returned. However, if you index the database using the NYSIIS algorithm and search using the NYSIIS algorithm again, the correct match will be returned because both "John Smith" and "Jon Smath" are indexed as "JANSNATH" by the algorithm. This option was developed to respond to limitations of Soundex; it handles some multicharacter n-grams and maintains relative vowel positioning, whereas Soundex does not.</p> <p>Note: This algorithm does not process non-alpha characters; records containing them will fail during processing.</p>
Phonix	It preprocesses name strings by applying more than 100 transformation rules to single characters or sequences of several characters. Nineteen of those rules are applied only if the characters are at the beginning of the string, while 12 of the rules are applied only if they are at the middle of the string, and 28 of the rules are applied only if they are at the end of the string. The transformed name string is encoded into a code that is comprised of a starting letter followed by three digits (removing zeros and duplicate numbers). This option was

developed to respond to the limitations of Soundex; it is more complex and, therefore, slower than Soundex.

Sonnex	It determines the similarity between two French-language strings based on the phonetic representation of their characters. It returns a Sonnex coded key of the selected fields.
Soundex	It determines the similarity between two strings based on a phonetic representation of their characters.
Syllable Alignment	It combines phonetic information with edit distance based calculations. Converts the strings to be compared into their corresponding sequences of syllables and calculates the number of edits required to convert one sequence of syllables to the other.

Similarity and Distance Measures

Edit Distance	It determines the similarity between two strings based on the number of deletions, insertions, or substitutions required to transform one string into another.
Euclidean Distance	<p>It provides a similarity measure between two strings using the vector space of combined terms as the dimensions. It also determines the greatest common divisor of two integers. It takes a pair of positive integers and forms a new pair that consists of the smaller number and the difference between the larger and smaller numbers. The process repeats until the numbers are equal. That number then is the greatest common divisor of the original pair.</p> <p>For example, 21 is the greatest common divisor of 252 and 105: $(252 = 12 \times 21; 105 = 5 \times 21)$; since $252 - 105 = (12 - 5) \times 21 = 147$, the GCD of 147 and 105 is also 21.</p>
Jaro-Winkler Distance	It determines the similarity between two strings based on the number of character replacements it takes to transform one string into another. This option was developed for short strings, such as personal names.
Keyboard Distance	<p>It determines the similarity between two strings based on the number of deletions, insertions, or substitutions required to transform one string to the other, weighted by the position of the keys on the keyboard.</p> <p>Click Edit to specify the type of keyboard you are using: QWERTY (U.S.), QWERTZ (Austria and Germany), or AZERTY (France).</p>
Kullback-Liebler Distance	It determines the similarity between two strings based on the differences between the distribution of words in the two strings.
NGram Distance	<p>It calculates in text or speech the probability of the next term based on the previous n terms, which can include phonemes, syllables, letters, words, or base pairs and can consist of any combination of letters.</p> <p>Click Edit to enter the size of the NGram; the default is 2.</p>

NGram Similarity It determines the similarity between two strings based on the length of the longest common subsequence of phonemes, syllables, letters, words, or base pairs.

Click **Edit** to specify these options:

- **Ngram size:** Enter the size of the NGram. The default is 2.
- **Drop Noise Characters:** Select to replace punctuation with space.
- **Drop Spaces:** Select to merge words.

Date Algorithms

Date It compares date fields regardless of the date format in the input records. Click **Edit** to specify these:

- **General Options—Require Month:** It prevents a date that consists only of a year from matching.
- **General Options—Require Day:** It prevents a date that consists only of a month and year from matching.
- **General Options—Match Transposed MM/DD:** Where month and day are provided in numeric format, it compares suspect month to candidate day and suspect day to candidate month as well as the standard comparison of the suspect month to candidate month and suspect day to candidate day.
- **General Options—Prefer DD/MM/YYYY format over MM/DD/YYYY:** It contributes to date parsing in cases where both month and day are provided in numeric format, and their identification can not be determined by context.

For example, given the numbers 5 and 13, the parser will automatically assign 5 to the month and 13 to the day because there are only 12 months in a year. However, given the numbers 5 and 12 (or any two numbers 12 and under), the parser will assume whichever number is first to be the month.

If you select this option, it ensures that the parser reads the first number as the day rather than the month.

- **Range Options—Overall:** It allows you to set the maximum number of days between matching dates. See the examples below.
 - If you enter an overall range of 35 days and your candidate date is December 31, 2000, a suspect date of February 5, 2001, would be a match, but a suspect date of February 6 would not.
 - If you enter an overall range of 1 day and your candidate date is January 2000, a suspect date of 1999 would be a match (comparing December 31, 1999), but a suspect date of January 2001 would not.
- **Range Options—Year:** It allows you to set the number of years between matching dates, independent of month, and day. See the examples below.

- If you enter a year range of 3 and your candidate date is January 31, 2000, a suspect date of January 31, 2003, would be a match, but a suspect date of February 2003 would not.
- If your candidate date is 2000, a suspect date of March 2003 would be a match because months are not in conflict, and it's within the three-year range.
- **Range Options—Month:** It allows you to set the number of months between matching dates, independent of year and day.

For example, if you enter a month range of 4 and your candidate date is January 1, 2000, a suspect date of May 2000 is a match because there is no day conflict and it's within the four-month range, but a suspect date of May 2, 2000, is not, because of the day's conflict.

- **Range Options—Day:** It allows you to set the number of days between matching dates, independent of year and month.

For example, if you enter a day range of 5 and your candidate date is January 1, 2000, a suspect date of January 2000 is a match because there is no day conflict but a suspect date of December 27, 1999, is not, because of the month's conflict.

Matching Algorithm-to-Scoring Method Matrix

The table below describes the logical relationship between the number of algorithms you can use based on the parent **Scoring Method** you select.

Scoring Method	Algorithms	
	Single	Multiple
Weighted Average	NA	Yes
Average	NA	Yes
Vector Summation	Yes	Yes
Maximum	NA	Yes
Minimum	NA	Yes

6. Click **Ok**.

Note:

- If you define n number of parent and child elements, use **Filter**  to selectively look for the elements.
- If you want to expand or collapse all the tree nodes, click the **Expand all** and **Collapse all**   buttons.

7. (Optional) If you want to store the match rule in the repository for re-use, click **Publish**.

The match rule is stored in the **Repository Rules** panel located on the left corner of the page.

8. On the top-right corner of the page, click **Apply**.

Your match rule applies to the **Intraflow Match** stage in the dataflow.

Evaluating a match rule

The **Evaluate** option allows you to test the results for your match rule. You can test any node you want in the match rule hierarchy. To evaluate, perform these steps:

1. In the match rule hierarchy, choose the node you want to test and click **Evaluate**.

You see the **Evaluate** page.

2. In the **Input** panel, enter the test data in these ways:

- To enter the test data manually:
 - a. Enter a suspect record under the **Suspect** column and up to ten candidates under **Candidate** columns.
 - b. Click **Export** to save the records to a file that you can import later instead of reentering the data manually.
- To import the test data from a file:
 - a. Click **Import**.
 - b. Select the file that contains the sample records. Delimited files can be comma, pipe, or tab-delimited and should have a header record with header fields that match the field names shown under **Candidates**. For example, a sample header record for household input would be `Name, AddressLine1, City, StateProvince`.

3. Select any of these method to evaluate your match rule:

- a. **Current Rule:** This runs the rule defined on the **Match Rule** panel. The results are displayed for one suspect and candidate pair at a time. To cycle through the results, select the arrows  . Scores for fields and algorithms are displayed in a tree format similar to the match rule control.

Note: If you make changes to the match rule and want to apply the changes to the stage's match rule, click **Save** on the top-right corner of the page.

- b. **All Algorithms:** This ignores the match rule and instead runs all algorithms against each field for suspect and candidate pairs. Results are displayed for one suspect and candidate pair at a time and can be cycled through using the arrows.

Your results are displayed in the **Scores** panel on the right corner of the screen with these color codes:

- **Green:** The rule resulted in a match
- **Red:** The rule did not result in a match
- **Gray:** The rule was ignored
- **Purple:** The results for individual algorithms within the rule

Output

For more information, see [Output](#) on page 292.

Best of Breed

The **Best of Breed** stage consolidates duplicate records by selecting the best data in a duplicate record collection and creating a new consolidated record using the best data.

This consolidated record is known as the *best of breed* record. You define the rules that select records to consolidate. When processing completes, the best of breed record is retained by the system.

This stage allows you to:

- **Create and publish a rule:** Define your rule in the **Rule Configuration** panel, and publish it to the repository for re-use at a later point in time, which allows you to save significant time. You can also use the **Create new rule** button at the left corner of the page to remove the current rule and define a new rule.

You can use **Template Rules**, which gives you the flexibility to choose a sophisticated set of predefined rules that saves time than defining a rule from scratch.

- Configure the options for the rule in the **Settings** panel. For more information, see [Configuring the Options](#) on page 411.
- **Configure access control for the rule:** It allows you to regulate who can view your rule and define the permissions for them. For more information, see [Configuring Access Control for Rule](#) on page 423.

Configuring the Options

On the **Best of Breed Options: Best of Breed** page, configure these options in the **Settings** panel:

Option	Description
Group by	<p>It specifies the field to use to create groups of records to merge into a single best of breed record, creating one best of breed record from each group. In cases where you have used a matching stage earlier in the dataflow, you should select the CollectionNumber field to use the collections created by the matching stage as the groups. However, if you want to group records by some other field, choose the field here.</p> <p>For example, if you want to merge all records that have the same value in the AccountNumber field into one best of breed record, you would select AccountNumber.</p>
Sort	<p>If you specify a field in the Group by field, check this box to sort the records by the value in the field you choose.</p>

Option	Description
Advanced	<p>Click this button to specify additional sort options. By default, the sort performance options specified in Management Console, which are the default performance options for your system, are in effect. If you want to override your system's default performance options, toggle the Override sort performance options to ON and specify the values you want in these fields:</p>
In memory record limit	<p>Specifies the maximum number of data rows a sorter will hold in memory before it starts paging to disk. By default, a sort of 10,000 records or less will be done in memory, and a sort of more than 10,000 records will be performed as a disk sort. The maximum limit is 100,000 records. An in-memory sort is much faster than a disk sort, so this value should be set high enough so that most of the sorts are in-memory sorts, and only large sets are written to disk.</p> <p>Note: Be careful in environments where jobs are running concurrently because increasing the In memory record limit setting increases the likelihood of running out of memory.</p>
Maximum number of temporary files	<p>Specifies the maximum number of temporary files that may be used by a sort process. Using a larger number of temporary files can result in better performance. However, the optimal number is highly dependent on the configuration of the server running Spectrum Technology Platform. You should experiment with different settings, observing the effect on the performance of using more or fewer temporary files. To calculate the approximate number of temporary files that may be needed, use this equation:</p>
	$\frac{(NumberOfRecords \times 2)}{InMemoryRecordLimit} = NumberOfTempFiles$
	<p>Note: The maximum number of temporary files cannot be more than 1,000.</p>
Compression	<p>Specifies that temporary files are compressed when they are written to disk.</p> <p>Note: The optimal sort of performance settings depends on your server's hardware configuration. You can use this equation as a general guideline to produce good sort performance: $(InMemoryRecordLimit \times MaxNumberOfTempFiles \div 2) \geq TotalNumberOfRecords$</p>
Keep original records	<p>Check this option to retain all records in the collection along with the best of breed record. Clear this option if you want only the best of breed record.</p>

Option	Description
Track master record	<p>Check this option to track the source of Best of Breed record. By enabling this option, you can select a unique input field from the drop-down to be used as a source record identifier. The selected values are then displayed in these output fields - AllSourceRecordsID and SelectedSourceRecordsID. For more information, see Output on page 424.</p> <p>Note: Ensure that you select a unique input field from the drop-down list.</p>

Defining Template Record Settings

In Best of Breed processing, the **Template Record** is the record in a collection that is used to create the best of breed record.

The template record is used as the starting point for constructing the best of breed record and is modified based on the **Best of Breed Settings** you define. The Best of Breed stage can select the template record automatically, or you can define rules for selecting the template record.

Example

Template record rules are written by specifying these options - **Field name**, **Field type**, **Operator**, **Value type**, and **Value**.

If you define the template record options as below:

Field name: MatchScore

Field type: Numeric

Operator: Equal

Value type: String

Value: 100

The template rule will select the record in the collection where the Match Score is equal to the value of 100.

To define a template record rule, go to **Template Record Settings** in the **Rule Configuration** panel and perform these steps:

1. Select the option **Define template record**.
2. Click the **Add rule**  button.
3. Define these options as listed in the table below, which are displayed on the **Rules** node:

Option	Description
Field name	It specifies the name of the dataflow field whose value you want to evaluate to determine if the record should be the template record.
Field type	Select one of these from the drop-down list to specify the type of data in the field: <ul style="list-style-type: none">• Non-Numeric: Select this option if the field contains non-numeric data, for example, string data.• Numeric: Choose this option if the field contains numeric data, for example, double or float.

Option	Description
Operator	<p>Select one of these from the drop-down list to specify the type of comparison you want to use to evaluate the field:</p> <ul style="list-style-type: none"> • Contains: Ignores the field if it contains blank data. • Equal: Scores the field as 0 if it contains blank data. • Greater Than: Scores the field as 100 if it contains blank data. • Greater Than Or Equal To: Scores the suspect and candidate fields as 100 if they both contain blank data; otherwise, scores the suspect and candidate fields as 0. • Highest: Compares the field's value for all the records group and determines which record has the highest value in the field. <p>For example, if the fields in the group contain 10, 20, 30, or 100, the record with the field value 100 would be selected.</p> <p>This operation only works on numeric fields. If multiple records are tied for the longest value, one record is selected.</p> <ul style="list-style-type: none"> • Is Empty: Determines if the field contains no value. • Is Not Empty: Determines if the field contains any value. • Less Than: Determines if the field value is less than the value specified. This operation only works on numeric fields. • Less Than Or Equal To: Determines if the field value is less than or equal to the value specified. This operation only works on numeric fields. • Longest: Compares the field's value for all the records group and determines which record has the longest (in bytes) value in the field. <p>For example, if the group contains the values "Mike" and "Michael", the record with the value "Michael" would be selected. If multiple records are tied for the longest value, one record is selected.</p> <ul style="list-style-type: none"> • Shortest: Compares the field's value for all the records group and determines which record has the shortest (in bytes) value in the field. <p>For example, if the group contains the values "Mike" and "Michael", the record with the value "Mike" would be selected. If multiple records are tied for the shortest value, one record is selected.</p> <ul style="list-style-type: none"> • Lowest: Compares the field's value for all the records group and determines which record has the lowest value in the field. <p>For example, if the fields in the group contain 10, 20, 30, or 100, the record with the field value 10 would be selected. This operation only works on numeric fields. If multiple records are tied for the longest value, one record is selected.</p> <ul style="list-style-type: none"> • Most Common: Determines if the field value contains the value that occurs most frequently in this field among the group's records. If you select this operator, you see the Multiple results option, which allows you to choose one of these: <ul style="list-style-type: none"> a. Consider None: If the field has the same number of occurrences for more than one value, no value is returned as most common. b. Consider All: If the field has the same number of occurrences for more than one value, all most common values are returned. • Not Equal: Determines if the field value is not the same as the value specified.

Option	Description
Value type	<p>Select one of these from the drop-down list to specify the value you want to compare to the field's value:</p> <p>Note: This option is not available if you select one of these Operators - Highest, Lowest, Longest, Shortest, Is Empty, Is Not Empty, or Most Common.</p> <ul style="list-style-type: none"> • Field: Select this option if you want to compare another dataflow field's value to the field. • String: Select this option if you want to compare the field to a specific value.
Value	<p>It specifies the value to compare to the field's value.</p> <ul style="list-style-type: none"> • If you selected Field in the Value type field, select a dataflow field. • If you selected String in the Value type field, enter the value you want to use in the comparison. <p>Note: This option is not available if you select one of these Operators - Highest, Lowest, Longest, Shortest, Is Empty, Is Not Empty, or Most Common.</p>

4. (Optional) If you want to define additional rules, click the **Add rule**  button.

Note: If you add additional rules, you will have to select a logical operator to use between each rule.

- If you want the new rule and the previous rule to both pass in order for it to be selected as the template record, select **And**.
- If you want either the previous rule or the new rule to pass in order for the record to be selected as the template record, select **Or**.

You have now configured rules to use to select the template record. You must now configure the best of breed settings to complete the configuration of the **Best of Breed** stage.

Defining Best of Breed Settings

The Best of Breed rules and actions work together to determine which fields from duplicate records in a collection to copy to the Best of Breed record. Rules test values in a record, and if the record passes the rules, the data is copied from the record to the template record. Actions define which data to copy, and which field in the template record should receive the data. After all the rules and actions are run, the template record will be the best of breed record.

Note: In **Best of Breed Settings, Rules and Actions** are grouped together into conditions, and you can have multiple conditions.

To define Best of Breed rules and actions, go to **Best of Breed Settings** in the **Rule Configuration** panel and perform these steps:

1. Click the **Add condition**  button.

A condition node is displayed to define the rules and actions.

2. Click the **Add rule**  button.
3. Define these options as listed in the table below, which are displayed on the **Rules** node:

Option	Description
Field name	It specifies the name of the dataflow field whose value you want to evaluate to determine if the condition is met and the associated actions should be taken.
Field type	Select one of these from the drop-down list to specify the type of data in the field: <ul style="list-style-type: none">• Non-Numeric: Select this option if the field contains non-numeric data, for example, string data.• Numeric: Choose this option if the field contains numeric data, for example, double or float.

Option	Description
--------	-------------

Operator	
-----------------	--

Option	Description
	<p>Select one of these from the drop-down list to specify the type of comparison you want to use to evaluate the field:</p> <ul style="list-style-type: none"> • Contains: Determines if the field contains the value specified, for example, "sailboat" contains the value "boat". • Equal: Determines if the field contains the exact value specified. • Greater Than: Determines if the field value is greater than the value specified. This operation only works on numeric fields. • Greater Than Or Equal To: Determines if the field value is greater than or equal to the value specified. <p style="text-align: center;">Note: This operation only works on numeric fields.</p>
	<ul style="list-style-type: none"> • Highest: Compares the field's value for all the records group and determines which record has the highest value in the field. <p>For example, if the fields in the group contain 10, 20, 30, or 100, the record with the field value 100 would be selected.</p> <p>This operation only works on numeric fields. If multiple records are tied for the longest value, one record is selected.</p>
	<ul style="list-style-type: none"> • Is Empty: Determines if the field contains no value. • Is Not Empty: Determines if the field contains any value. • Less Than: Determines if the field value is less than the value specified. This operation only works on numeric fields. • Less Than Or Equal To: Determines if the field value is less than or equal to the value specified. This operation only works on numeric fields. • Longest: Compares the field's value for all the records group and determines which record has the longest (in bytes) value in the field. <p>For example, if the group contains the values "Mike" and "Michael", the record with the value "Michael" would be selected. If multiple records are tied for the longest value, one record is selected.</p>
	<ul style="list-style-type: none"> • Shortest: Compares the field's value for all the records group and determines which record has the shortest (in bytes) value in the field. <p>For example, if the group contains the values "Mike" and "Michael", the record with the value "Mike" would be selected. If multiple records are tied for the shortest value, one record is selected.</p>
	<ul style="list-style-type: none"> • Lowest: Compares the field's value for all the records group and determines which record has the lowest value in the field. <p>For example, if the fields in the group contain 10, 20, 30, or 100, the record with the field value 10 would be selected. This operation only works on numeric fields. If multiple records are tied for the longest value, one record is selected.</p> <ul style="list-style-type: none"> • Most Common: Determines if the field value contains the value that occurs most frequently in this field among the group's records. If you select this operator, you see the Multiple results option, which allows you to choose one of these: <ul style="list-style-type: none"> a. Consider None: If the field has the same number of occurrences for more than one value, no value is returned as most common. b. Consider All: If the field has the same number of occurrences for more than

Option	Description
	<p>one value, all most common values are returned.</p> <ul style="list-style-type: none"> • Not Equal: Determines if the field value is not the same as the value specified.
Value type	<p>Select one of these from the drop-down list to specify the value you want to compare to the field's value:</p> <p>Note: This option is not available if you select one of these Operators - Highest, Lowest, Longest, Shortest, Is Empty, Is Not Empty, or Most Common.</p> <ul style="list-style-type: none"> • Field: Select this option if you want to compare another dataflow field's value to the field. • String: Select this option if you want to compare the field to a specific value.
Value	<p>It specifies the value to compare to the field's value.</p> <ul style="list-style-type: none"> • If you selected Field in the Value type field, select a dataflow field. • If you selected String in the Value type field, enter the value you want to use in the comparison. <p>Note: This option is not available if you select one of these Operators - Highest, Lowest, Longest, Shortest, Is Empty, Is Not Empty, or Most Common.</p>

- (Optional) If you want to define additional rules, click the **Add rule** button.

Note: If you add additional rules, you will have to select a logical operator to use between each rule.

- If you want the new rule and the previous rule to both pass in order for the condition to be met and the associated actions taken, select **And**.
- If you want either the previous rule or the new rule to pass in order for the condition to be met. select **Or**.

- Click the **Add action** button.
- Define these options as listed in the table below, which are displayed on the **Actions** node:

Option	Description
Source type	<p>Select one of these from the drop-down list to specify the type of data to copy to the Best of Breed record:</p> <ul style="list-style-type: none"> • Field: Select this option if you want to copy a value from a field to the Best of Breed record. • String: Select this option if you want to copy a constant value to the Best of Breed record.
Source data	<p>It specifies the data to copy to the best of breed record.</p> <ul style="list-style-type: none"> • If the Source type is Field, select the field whose value you want to copy to the destination field. • If the Source type is String, specify a constant value to copy to the destination field.
Destination	<p>It specifies the field in the Best of Breed record to which you want to copy the data defined in the Source data field.</p>
Accumulate source data	<p>If the data in the Source data field is numeric, you can enable this option to combine the source data for all duplicate records and put the total value in the Best of Breed record.</p> <p>For example, if there were three duplicate records in the group and they contained these values in the Deposits field:</p> <p>100.00 20.00 5.00</p> <p>Then all three values would be combined, and the total value, 125.00, would be put in the Best of Breed record's Deposits field.</p>

7. (Optional) If you want to define additional actions for the same condition , repeat steps **5** on page 421 and **6** on page 421.
8. (Optional) If you want to add another condition, click the **Add condition**  button.
9. (Optional) If you want to store the rule in the repository for re-use:
 - a. Click **Publish**.
You see a **Publish Best of Breed Rule** panel on the right corner of the page.
 - b. Enter a name for your rule in the **Rule name** field.
 - c. Click **OK**.

The rule is stored and appears in **Repository Rules** located on the left corner of the page, which allows you to save significant time by re-using the rules at a later point in time. If you modify the repository rule, by default, **Custom** is added as a prefix to the Rule name, for example, Custom - xyz.

Note: You can regulate the access to your rule if you do not want any other user to re-use them. For more information, see [Configuring Access Control for Rule](#) on page 423.

- On the top-right corner of the page, click **Apply**.

Your rule applies to the Best of Breed stage in the dataflow.

Example for Best of Breed Rule and Action

This Best of Breed rule selects the record where the Match Score is equal to the value of 100. The Account Number data that corresponds to the selected field is then copied to the AccountNumber field on the Best of Breed record.

Rule

Field name: MatchScore

Field type: Numeric

Operator: Equal

Value type: String

Value: 100

Action

Source type: Field

Source data: AccountNumber

Destination: AccountNumber

Configuring Access Control for Rule

Access control to a rule minimizes the risk of unauthorized usage by a user. It allows you to regulate who can view your rule and define the permissions for them.

To configure access control for your rule, perform these steps:

- Go to **Management Console > System > Security**.
- Click the **Access Control** tab.
- Click the **Add access control**  button.
- Do one of these:
 - If you want to specify access controls for a role, select **Role** from the drop-down field. The access control permissions you specify will affect all users who have the role you choose.
 - If you want to specify access controls for a single user, select **User** from the drop-down field. The access control permissions you specify will only affect the user you choose.
- Click the **Add entity override**  button.
You see an **Add Entity** workspace.
- In the **Add Entity** workspace, select **AMM.Best of Breed Rules Management** from the drop-down field and do these:

- a. Drag the rules you see on the left side to the **Selected Entities** space on the right side.
- b. Click **Add**.

The rules you add are displayed in the **AMM.Best of Breed Rules Management** panel at the bottom on the **Add Access Control** page.

7. Specify the **View**, **Modify**, and **Delete** permissions that you want to grant for each rule. You can select one of these options for each permission:

Option	Description
	The permission is inherited from the role.
	The permission is inherited from the role and cannot be overridden.
	The permission is granted, overriding the permission specified in the user or role.
	The permission is denied, overriding the permission specified in the user or role.

8. On the top-right corner of the page, select **Save**. Your configuraton applies to the rules.

Output

Option	Format	Description
CollectionRecordType	String	<p>It identifies the template and Best of Breed records in a collection of duplicate records. The possible values are:</p> <ul style="list-style-type: none"> • Primary: The record is the selected template record in a collection. • Secondary: The record is not the selected template record in a collection. • BestOfBreed: The record is the newly created Best of Breed record in the collection. <p>Note: The Primary and Secondary values are generated only when you select the Define template record option in Template Record Settings.</p>
AllSourceRecordsID		It displays the value of all the records for the Best of Breed record, which are in the same group.

Option	Format	Description
SelectedSourceRecordsID		It displays the value of the selected records for the Best of Breed record.

11 - Performance Tuning

In this section

Performance Tuning Recommendations.....427



Performance Tuning Recommendations

The sections below combine many tips that might help you improve the performance of various data quality stages. The recommendations listed will help you achieve performance for **Advanced Matching**, **Data Normalization**, and **Universal Name** stages.

The recommendations below are applicable to all stages:

- **Number of records:** Analyzing and filtering the records before sending them for processing improves the performance, as an increase in the number of records leads to a proportional increase in time.
- **Cluster:** The performance increases when the processing is performed in clustered mode.

Performance recommendations for Advanced Matching stages

Stage	Performance recommendations
Intraflow Match	<ul style="list-style-type: none"> • Group size: Minimizing the group size or match key to the optimum level leads to better performance, as, given the same number of records, an increase in group size leads to increased processing time. • Express key: Using an express key during matching improves the performance, but it should be evaluated to ensure that the express key is an accurate candidate for the express match. • Input data: Receiving input data sorted by match key improves the performance. • Match rule: Optimizing the match rule increases the performance, as a complex match rule degrades the performance.
Interflow Match	<ul style="list-style-type: none"> • Group size: Minimizing the group size or match key to the optimum level leads to better performance, as, given the same number of records, an increase in group size leads to increased processing time. • Match rule and key: Optimizing the match rule and match keys well, as it plays a crucial role in increasing the performance.
Transactional Match	<p>This stage is fast as compared to the Intraflow Match stage; however, it is recommended to choose the relevant stage according to your requirements as both serve two different purposes. Moreover, unlike Intraflow Match, the group size does not affect the performance as the suspect is matched only once with the candidate.</p>

Stage	Performance recommendations
Best of Breed Duplicate Synchronization Filter	Condition: Minimizing the conditions improves the performance, as an increased number of conditions leads to more processing time.
Match Key Generator	Runtime: An increase in runtime instances increases performance.
Candidate Finder	<p>The recommendations below improves the performance for Search Index.</p> <ul style="list-style-type: none"> Runtime: The search operation's performance increases when the runtime instances increases. The machine's configuration decides on how many runtime instances can be used. <p>For example, we observed performance improvement with an increase in the runtime instances.</p> Fields: Create and search operations' performance decreases with more number of fields in the index. However, update operation performance remains almost the same irrespective of the number of fields in the index. <p>For example, we observed performance degradation in the case of a search when the number of fields was increased.</p> Batch size: Performance varies when the batch size changes. The optimal number of batch size for the machine is figured out with different batch values. It depends on the machine's memory and CPU resources. Shards: Update operation improves to some level when the number of shards increases. Search performance degrades with more number of shards. <p>For example, we observed that the update gets relatively faster when the shards increase, whereas the search performance degrades.</p> Candidate Finder (CF) conditions: The search takes time when the number of conditions in the candidate finder stage increases. <p>For example, we observed performance degradation when the CF query increased.</p> Analyzer: Search using a keyword analyzer is much faster than a standard analyzer. <p>For example, we observed a performance improvement when the analyzer changed from standard to keyword.</p> <p>Note: Search index performance depends on various factors, and the points mentioned above are just an indication of how the performance can vary based on the configuration applied. It is essential to understand the end-user scenario, which is the key driver for deciding on the choice of hardware, index settings, cluster setup, and other configuration parameters to achieve optimal performance.</p>

Performance recommendation for Data Normalization stages

Stage	Performance recommendation
Table Lookup Advanced Transformer Open Parser	Runtime: An increase in runtime instances increases performance.

Performance recommendation for Universal Name stage

Stage	Performance recommendation
Name Parser	Runtime: An increase in runtime instances increases performance.

12 - ISO Country Codes and Module Support

In this section

ISO Country Codes and Module Support.....431



ISO Country Codes and Module Support

This table lists the ISO codes for each country as well as the modules that support addressing, geocoding, and routing for each country.

Note that the Spectrum Enterprise Geocoding includes databases for Africa (30 countries), Middle East (8 countries) and Latin America (20 countries). These databases cover the smaller countries in those regions that do not have their own country-specific geocoding databases. The Supported Modules column indicates which countries are covered by these Africa, Middle East, and Latin America databases.

Also, the Geocode Address World database provides geographic and limited postal geocoding (but not street-level geocoding) for all countries.

ISO Country Name	ISO 3166-1 Alpha-2	ISO 3166-1 Alpha-3	Supported Modules
Afghanistan	AF	AFG	Spectrum Universal Address
Aland Islands	AX	ALA	Spectrum Universal Address
Albania	AL or SQ (Routing)	ALB	Spectrum Universal Address Spectrum Enterprise Geocoding Spectrum Spatial Routing
Algeria	DZ	DZA	Spectrum Enterprise Geocoding (Africa) Spectrum Universal Address
American Samoa	AS	ASM	Spectrum Universal Address
Andorra	AD	AND	Spectrum Enterprise Geocoding. (Andorra is covered by the Spain geocoder) Spectrum Universal Address

ISO Country Name	ISO 3166-1 Alpha-2	ISO 3166-1 Alpha-3	Supported Modules
Angola	AO	AGO	Spectrum Enterprise Geocoding (Africa) Spectrum Universal Address
Anguilla	AI	AIA	Spectrum Universal Address
Antarctica	AQ	ATA	Spectrum Universal Address
Antigua And Barbuda	AG	ATG	Spectrum Universal Address
Argentina	AR	ARG	Spectrum Enterprise Geocoding Spectrum Universal Address Spectrum Spatial Routing
Armenia	AM	ARM	Spectrum Universal Address
Aruba	AW	ABW	Spectrum Enterprise Geocoding (Latin America) Spectrum Universal Address
Australia	AU	AUS	Spectrum Enterprise Geocoding Spectrum Spatial Routing Spectrum Universal Address
Austria	AT	AUT	Spectrum Enterprise Geocoding Spectrum Spatial Routing Spectrum Universal Address
Azerbaijan	AZ	AZE	Spectrum Universal Address

ISO Country Name	ISO 3166-1 Alpha-2	ISO 3166-1 Alpha-3	Supported Modules
Bahamas	BS	BHS	Spectrum Enterprise Geocoding Spectrum Universal Address Spectrum Spatial Routing
Bahrain	BH	BHR	Spectrum Enterprise Geocoding (Middle East) Spectrum Universal Address
Bangladesh	BD	BGD	Spectrum Universal Address
Barbados	BB	BRB	Spectrum Enterprise Geocoding (Latin America) Spectrum Universal Address
Belarus	BY	BLR	Spectrum Universal Address Spectrum Spatial Routing
Belgium	BE	BEL	Spectrum Enterprise Geocoding Spectrum Spatial Routing Spectrum Universal Address
Belize	BZ	BLZ	Spectrum Enterprise Geocoding (Latin America) Spectrum Universal Address
Benin	BJ	BEN	Spectrum Enterprise Geocoding (Africa) Spectrum Universal Address
Bermuda	BM	BMU	Spectrum Universal Address Spectrum Spatial Routing
Bhutan	BT	BTN	Spectrum Universal Address

ISO Country Name	ISO 3166-1 Alpha-2	ISO 3166-1 Alpha-3	Supported Modules
Bolivia	BO	BOL	Spectrum Enterprise Geocoding (Latin America) Spectrum Universal Address
Bonaire, Saint Eustatius And Saba	BQ	BES	Spectrum Universal Address
Bosnia And Herzegovina	BA	BIH	Spectrum Enterprise Geocoding Spectrum Universal Address Spectrum Spatial Routing Spectrum Enterprise Geocoding
Botswana	BW	BWA	Spectrum Enterprise Geocoding (Africa) Spectrum Universal Address
Bouvet Island	BV	BVT	Spectrum Universal Address
Brazil	BR	BRA	Spectrum Enterprise Geocoding Spectrum Spatial Routing Spectrum Universal Address
British Indian Ocean Territory	IO	IOT	Spectrum Universal Address
Brunei Darussalam	BN	BRN	Spectrum Enterprise Geocoding Spectrum Universal Address
Bulgaria	BG	BGR	Spectrum Enterprise Geocoding Spectrum Universal Address
Burkina Faso	BF	BFA	Spectrum Enterprise Geocoding (Africa) Spectrum Universal Address

ISO Country Name	ISO 3166-1 Alpha-2	ISO 3166-1 Alpha-3	Supported Modules
Burundi	BI	BDI	Spectrum Enterprise Geocoding (Africa) Spectrum Universal Address
Cambodia	KH	KHM	Spectrum Universal Address
Cameroon	CM	CMR	Spectrum Enterprise Geocoding (Africa) Spectrum Universal Address
Canada	CA	CAN	Spectrum Enterprise Geocoding Spectrum Spatial Routing Spectrum Universal Address
Cape Verde	CV	CPV	Spectrum Universal Address
Cayman Islands	KY	CYM	Spectrum Universal Address
Central African Republic	CF	CAF	Spectrum Universal Address
Chad	TD	TCD	Spectrum Universal Address
Chile	CL	CHL	Spectrum Enterprise Geocoding Spectrum Universal Address Spectrum Spatial Routing
China	CN or zh_CN (Routing)	CHN	Spectrum Enterprise Geocoding Spectrum Universal Address Spectrum Spatial Routing
Christmas Island	CX	CXR	Spectrum Universal Address

ISO Country Name	ISO 3166-1 Alpha-2	ISO 3166-1 Alpha-3	Supported Modules
Cocos (Keeling) Islands	CC	CCK	Spectrum Universal Address
Colombia	CO	COL	Spectrum Enterprise Geocoding Spectrum Universal Address
Comoros	KM	COM	Spectrum Universal Address
Congo, Republic Of The	CG	COG	Spectrum Enterprise Geocoding (Africa) Spectrum Universal Address
Congo, The Democratic Republic Of The	CD	COD	Spectrum Enterprise Geocoding (Africa) Spectrum Universal Address Spectrum Spatial Routing
Cook Islands	CK	COK	Spectrum Universal Address
Costa Rica	CR	CRI	Spectrum Enterprise Geocoding (Latin America) Spectrum Universal Address
Côte d'Ivoire	CI	CIV	Spectrum Universal Address
Croatia	HR	HRV	Spectrum Enterprise Geocoding Spectrum Spatial Routing Spectrum Universal Address
Cuba	CU	CUB	Spectrum Enterprise Geocoding (Latin America) Spectrum Spatial Routing Spectrum Universal Address
Curacao	CW	CUW	Spectrum Universal Address

ISO Country Name	ISO 3166-1 Alpha-2	ISO 3166-1 Alpha-3	Supported Modules
Cyprus	CY	CYP	Spectrum Enterprise Geocoding Spectrum Universal Address
Czech Republic	CZ or CS (Routing)	CZE	Spectrum Enterprise Geocoding Spectrum Universal Address Spectrum Spatial Routing
Denmark	DK	DNK	Spectrum Enterprise Geocoding Spectrum Spatial Routing Spectrum Universal Address
Djibouti	DJ	DJI	Spectrum Universal Address
Dominica	DM	DMA	Spectrum Universal Address
Dominican Republic	DO	DOM	Spectrum Enterprise Geocoding (Latin America) Spectrum Universal Address
Ecuador	EC	ECU	Spectrum Enterprise Geocoding (Latin America) Spectrum Universal Address
Egypt	EG	EGY	Spectrum Enterprise Geocoding (Middle East) Spectrum Universal Address
El Salvador	SV	SLV	Spectrum Enterprise Geocoding (Latin America) Spectrum Universal Address
Equatorial Guinea	GQ	GNQ	Spectrum Universal Address
Eritrea	ER	ERI	Spectrum Universal Address

ISO Country Name	ISO 3166-1 Alpha-2	ISO 3166-1 Alpha-3	Supported Modules
Estonia	EE	EST	Spectrum Enterprise Geocoding Spectrum Spatial Routing Spectrum Universal Addressing
Ethiopia	ET	ETH	Spectrum Universal Address
Falkland Islands (Malvinas)	FK	FLK	Spectrum Universal Address
Faroe Islands	FO	FRO	Spectrum Universal Address
Fiji	FJ	FJI	Spectrum Universal Address
Finland	FI	FIN	Spectrum Enterprise Geocoding Spectrum Spatial Routing Spectrum Universal Address
France	FR	FRA	Spectrum Enterprise Geocoding Spectrum Spatial Routing Spectrum Universal Address
French Guiana	GF	GUF	Spectrum Enterprise Geocoding (<i>French Guiana is covered by the France geocoder.</i>) Spectrum Universal Address
French Polynesia	PF	PYF	Spectrum Universal Address
French Southern Territories	TF	ATF	Spectrum Universal Address
Gabon	GA	GAB	Spectrum Enterprise Geocoding (Africa) Spectrum Universal Address

ISO Country Name	ISO 3166-1 Alpha-2	ISO 3166-1 Alpha-3	Supported Modules
Gambia	GM	GMB	Spectrum Universal Address
Georgia	GE	GEO	Spectrum Universal Address
Germany	DE	DEU	Spectrum Enterprise Geocoding Spectrum Spatial Routing Spectrum Universal Address
Ghana	GH	GHA	Spectrum Enterprise Geocoding (Africa) Spectrum Universal Address Spectrum Spatial Routing
Gibraltar	GI	GIB	Spectrum Enterprise Geocoding (<i>Gibraltar is covered by the Spain geocoder.</i>) Spectrum Universal Address
Greece	GR	GRC	Spectrum Enterprise Geocoding Spectrum Universal Address
Greenland	GL	GRL	Spectrum Universal Address
Grenada	GD	GRD	Spectrum Universal Address
Guadeloupe	GP	GLP	Spectrum Enterprise Geocoding (<i>Guadeloupe is covered by the France geocoder.</i>) Spectrum Universal Address
Guam	GU	GUM	Spectrum Universal Address

ISO Country Name	ISO 3166-1 Alpha-2	ISO 3166-1 Alpha-3	Supported Modules
Guatemala	GT	GTM	Spectrum Enterprise Geocoding (Latin America) Spectrum Universal Address
Guernsey	GG	GGY	Spectrum Universal Address
Guinea	GN	GIN	Spectrum Universal Address
Guinea-Bissau	GW	GNB	Spectrum Universal Address
Guyana	GY	GUY	Spectrum Enterprise Geocoding (Latin America) Spectrum Universal Address
Haiti	HT	HTI	Spectrum Universal Address
Heard Island and McDonald Islands	HM	HMD	Spectrum Universal Address
Holy See (Vatican City State)	VA	VAT	Spectrum Enterprise Geocoding (<i>The Vatican is covered by the Italy geocoder.</i>) Spectrum Universal Address
Honduras	HN	HND	Spectrum Enterprise Geocoding (Latin America) Spectrum Universal Address
Hong Kong	HK	HKG	Spectrum Enterprise Geocoding Spectrum Universal Address
Hungary	HU	HUN	Spectrum Enterprise Geocoding Spectrum Spatial Routing Spectrum Universal Address

ISO Country Name	ISO 3166-1 Alpha-2	ISO 3166-1 Alpha-3	Supported Modules
Iceland	IS	ISL	Spectrum Enterprise Geocoding Spectrum Universal Address
India	IN	IND	Spectrum Enterprise Geocoding Spectrum Universal Address
Indonesia	ID	IDN	Spectrum Enterprise Geocoding Spectrum Universal Address
Iran, Islamic Republic Of	IR	IRN	Spectrum Universal Address
Iraq	IQ	IRQ	Spectrum Enterprise Geocoding (Middle East) Spectrum Universal Address
Ireland	IE	IRL	Spectrum Enterprise Geocoding Spectrum Spatial Routing Spectrum Universal Address
Isle Of Man	IM	IMN	Spectrum Universal Address
Israel	IL	ISR	Spectrum Enterprise Geocoding Spectrum Universal Address Spectrum Spatial Routing
Italy	IT	ITA	Spectrum Enterprise Geocoding Spectrum Spatial Routing Spectrum Universal Address
Jamaica	JM	JAM	Spectrum Enterprise Geocoding (Latin America) Spectrum Universal Address

ISO Country Name	ISO 3166-1 Alpha-2	ISO 3166-1 Alpha-3	Supported Modules
Japan	JP	JPN	Spectrum Enterprise Geocoding Spectrum Spatial Routing Spectrum Universal Address
Jersey	JE	JEY	Spectrum Universal Address
Jordan	JO	JOR	Spectrum Universal Address Spectrum Enterprise Geocoding (Middle East) Spectrum Spatial Routing
Kazakhstan	KZ	KAZ	Spectrum Universal Address
Kenya	KE	KEN	Spectrum Enterprise Geocoding (Africa) Spectrum Universal Address Spectrum Spatial Routing
Kiribati	KI	KIR	Spectrum Universal Address
Korea, Democratic People's Republic Of	KP	PRK	Spectrum Universal Address
Korea, Republic Of	KR	KOR	Spectrum Enterprise Geocoding Spectrum Universal Address
Kosovo	Xk	XXK	Spectrum Enterprise Geocoding Spectrum Universal Address
Kuwait	KW	KWT	Spectrum Enterprise Geocoding (Middle East) Spectrum Universal Address

ISO Country Name	ISO 3166-1 Alpha-2	ISO 3166-1 Alpha-3	Supported Modules
Kyrgyzstan	KG	KGZ	Spectrum Universal Address
Lao People's Democratic Republic	LA	LAO	Spectrum Universal Address
Latvia	LV	LVA	Spectrum Enterprise Geocoding Spectrum Spatial Routing Spectrum Universal Address
Lebanon	LB	LBN	Spectrum Enterprise Geocoding (Middle East) Spectrum Universal Address
Lesotho	LS	LSO	Spectrum Enterprise Geocoding (Africa) Spectrum Universal Address Spectrum Spatial Routing
Liberia	LR	LBR	Spectrum Universal Address
Libyan Arab Jamahiriya	LY	LBY	Spectrum Universal Address
Liechtenstein	LI	LIE	Spectrum Enterprise Geocoding (<i>Liechtenstein is covered by the Switzerland geocoder.</i>) Spectrum Spatial Routing Spectrum Universal Address
Lithuania	LT	LTU	Spectrum Enterprise Geocoding Spectrum Spatial Routing Spectrum Universal Address
Luxembourg	LU	LUX	Spectrum Enterprise Geocoding (<i>Luxembourg is covered by the Belgium geocoder.</i>) Spectrum Spatial Routing Spectrum Universal Address

ISO Country Name	ISO 3166-1 Alpha-2	ISO 3166-1 Alpha-3	Supported Modules
Macao	MO	MAC	Spectrum Enterprise Geocoding Spectrum Universal Address
Macedonia, Former Yugoslav Republic Of	MK	MKD	Spectrum Enterprise Geocoding Spectrum Universal Address
Madagascar	MG	MDG	Spectrum Universal Address
Malawi	MW	MWI	Spectrum Enterprise Geocoding (Africa) Spectrum Universal Address
Malaysia	MY	MYS	Spectrum Enterprise Geocoding Spectrum Spatial Routing Spectrum Universal Address
Maldives	MV	MDV	Spectrum Universal Address
Mali	ML	MLI	Spectrum Enterprise Geocoding (Africa) Spectrum Universal Address
Malta	ML	MLT	Spectrum Enterprise Geocoding Spectrum Universal Address
Marshall Islands	MH	MHL	Spectrum Universal Address
Martinique	MQ	MTQ	Spectrum Enterprise Geocoding (<i>Martinique is covered by the France geocoder.</i>) Spectrum Universal Address
Mauritania	MR	MRT	Spectrum Enterprise Geocoding (Africa) Spectrum Universal Address

ISO Country Name	ISO 3166-1 Alpha-2	ISO 3166-1 Alpha-3	Supported Modules
Mauritius	MU	MUS	Spectrum Enterprise Geocoding (Africa) Spectrum Universal Address
Mayotte	YT	MYT	Spectrum Enterprise Geocoding (<i>Mayotte is covered by the France geocoder.</i>) Spectrum Universal Address
Mexico	MX	MEX	Spectrum Enterprise Geocoding Spectrum Universal Address
Micronesia, Federated States Of	FM	FSM	Spectrum Universal Address
Moldova, Republic Of	MD	MDA	Spectrum Universal Address Spectrum Spatial Routing
Monaco	MC	MCO	Spectrum Enterprise Geocoding (<i>Monaco is covered by the France geocoder.</i>) Spectrum Universal Address
Mongolia	MN	MNG	Spectrum Universal Address
Montenegro	ME	MNE	Spectrum Enterprise Geocoding Spectrum Universal Address
Montserrat	MS	MSR	Spectrum Universal Address
Morocco	MA	MAR	Spectrum Enterprise Geocoding (Africa) Spectrum Universal Address

ISO Country Name	ISO 3166-1 Alpha-2	ISO 3166-1 Alpha-3	Supported Modules
Mozambique	MZ	MOZ	Spectrum Enterprise Geocoding (Africa) Spectrum Universal Address Spectrum Spatial Routing
Myanmar	MM	MMR	Spectrum Universal Address
Namibia	NA	NAM	Spectrum Enterprise Geocoding (Africa) Spectrum Universal Address
Nauru	NR	NRU	Spectrum Universal Address
Nepal	NP	NPL	Spectrum Universal Address
Netherlands	NL	NLD	Spectrum Enterprise Geocoding Spectrum Spatial Routing Spectrum Universal Address
New Caledonia	NC	NCL	Spectrum Universal Address
New Zealand	NZ	NZL	Spectrum Enterprise Geocoding Spectrum Spatial Routing Spectrum Universal Address
Nicaragua	NI	NIC	Spectrum Enterprise Geocoding (Latin America) Spectrum Universal Address
Niger	NE	NER	Spectrum Enterprise Geocoding (Africa) Spectrum Universal Address

ISO Country Name	ISO 3166-1 Alpha-2	ISO 3166-1 Alpha-3	Supported Modules
Nigeria	NG	NGA	Spectrum Enterprise Geocoding (Africa) Spectrum Universal Address
Niue	NU	NIU	Spectrum Universal Address
Norfolk Island	NF	NFK	Spectrum Universal Address
Northern Mariana Islands	MP	MNP	Spectrum Universal Address
Norway	NO	NOR	Spectrum Enterprise Geocoding Spectrum Spatial Routing Spectrum Universal Address
Oman	OM	OMN	Spectrum Enterprise Geocoding (Middle East) Spectrum Universal Address
Pakistan	PK	PAK	Spectrum Universal Address
Palau	PW	PLW	Spectrum Universal Address
Palestinian Territory, Occupied	PS	PSE	Spectrum Universal Address
Panama	PA	PAN	Spectrum Enterprise Geocoding (Latin America) Spectrum Universal Address
Papua New Guinea	PG	PNG	Spectrum Universal Address

ISO Country Name	ISO 3166-1 Alpha-2	ISO 3166-1 Alpha-3	Supported Modules
Paraguay	PY	PRY	Spectrum Enterprise Geocoding (Latin America) Spectrum Universal Address
Peru	PE	PER	Spectrum Enterprise Geocoding (Latin America) Spectrum Universal Address
Philippines	PH	PHL	Spectrum Enterprise Geocoding Spectrum Universal Address Spectrum Spatial Routing
Pitcairn	PN	PCN	Spectrum Universal Address
Poland	PL	POL	Spectrum Enterprise Geocoding Spectrum Spatial Routing Spectrum Universal Address
Portugal	PT	PRT	Spectrum Enterprise Geocoding Spectrum Spatial Routing Spectrum Universal Address
Puerto Rico	PR	PRI	Spectrum Universal Address
Qatar	QA	QAT	Spectrum Enterprise Geocoding (Middle East) Spectrum Universal Address
Reunion	RE	REU	Spectrum Enterprise Geocoding (<i>Reunion is covered by the France geocoder.</i>) Spectrum Universal Address
Romania	RO	ROU	Spectrum Enterprise Geocoding Spectrum Spatial Routing Spectrum Universal Address

ISO Country Name	ISO 3166-1 Alpha-2	ISO 3166-1 Alpha-3	Supported Modules
Russian Federation	RU	RUS	Spectrum Enterprise Geocoding Spectrum Universal Address
Rwanda	RW	RWA	Spectrum Enterprise Geocoding (Africa) Spectrum Universal Address
Saint Barthelemy	BL	BLM	Spectrum Universal Address
Saint Helena, Ascension and Tristan Da Cunha	SH	SHE	Spectrum Universal Address
Saint Kitts and Nevis	KN	KNA	Spectrum Enterprise Geocoding (Latin America) Spectrum Universal Address
Saint Lucia	LC	LCA	Spectrum Universal Address
Saint Martin (French Part)	MF	MAF	Spectrum Universal Address
Saint Pierre and Miquelon	PM	SPM	Spectrum Universal Address
Saint Vincent and the Grenadines	VC	VCT	Spectrum Universal Address
Samoa	WS	WSM	Spectrum Universal Address
San Marino	SM	SMR	Spectrum Enterprise Geocoding (<i>San Marino is covered by the Italy geocoder.</i>) Spectrum Universal Address
Sao Tome and Principe	ST	STP	Spectrum Universal Address

ISO Country Name	ISO 3166-1 Alpha-2	ISO 3166-1 Alpha-3	Supported Modules
Saudi Arabia	SA	SAU	Spectrum Enterprise Geocoding (Middle East) Spectrum Universal Address
Senegal	SN	SEN	Spectrum Enterprise Geocoding (Africa) Spectrum Universal Address
Serbia	RS	SRB	Spectrum Enterprise Geocoding Spectrum Universal Address
Seychelles	SC	SYC	Spectrum Universal Address
Sierra Leone	SL	SLE	Spectrum Universal Address
Singapore	SG	SGP	Spectrum Enterprise Geocoding Spectrum Spatial Routing Spectrum Universal Address
Sint Maarten (Dutch Part)	SX	SXM	Spectrum Universal Address
Slovakia	SK	SVK	Spectrum Enterprise Geocoding Spectrum Spatial Routing Spectrum Universal Address
Slovenia	SI	SVN	Spectrum Enterprise Geocoding Spectrum Spatial Routing Spectrum Universal Address
Solomon Islands	SB	SLB	Spectrum Universal Address
Somalia	SO	SOM	Spectrum Universal Address

ISO Country Name	ISO 3166-1 Alpha-2	ISO 3166-1 Alpha-3	Supported Modules
South Africa	ZA	ZAF	Spectrum Enterprise Geocoding Spectrum Universal Address
South Georgia And The South Sandwich Islands	GS	SGS	Spectrum Enterprise Geocoding Spectrum Universal Address
South Sudan	SS	SSD	Spectrum Universal Address
Spain	ES	ESP	Spectrum Enterprise Geocoding Spectrum Spatial Routing Spectrum Universal Address
Sri Lanka	LK	LKA	Spectrum Universal Address
Sudan	SD	SDN	Spectrum Universal Address
Suriname	SR	SUR	Spectrum Enterprise Geocoding (Latin America) Spectrum Universal Address
Svalbard And Jan Mayen	SJ	SJM	Spectrum Universal Address
Swaziland	SZ	SWZ	Spectrum Enterprise Geocoding (Africa) Spectrum Universal Address
Sweden	SE	SWE	Spectrum Enterprise Geocoding Spectrum Spatial Routing Spectrum Universal Address
Switzerland	CH	CHE	Spectrum Enterprise Geocoding Spectrum Spatial Routing Spectrum Universal Address

ISO Country Name	ISO 3166-1 Alpha-2	ISO 3166-1 Alpha-3	Supported Modules
Syrian Arab Republic	SY	SYR	Spectrum Universal Address
Taiwan, Province of China	TW or zh_TW (Routing)	TWN	Spectrum Universal Address Spectrum Spatial Routing
Tajikistan	TJ	TJK	Spectrum Universal Address
Tanzania, United Republic Of	TZ	TZA	Spectrum Enterprise Geocoding (Africa) Spectrum Universal Address Spectrum Spatial Routing
Thailand	TH	THA	Spectrum Enterprise Geocoding Spectrum Spatial Routing Spectrum Universal Address
Timor-Leste	TL	TLS	Spectrum Universal Address
Togo	TG	TGO	Spectrum Enterprise Geocoding (Africa) Spectrum Universal Address
Tokelau	TK	TKL	Spectrum Universal Address
Tonga	TO	TON	Spectrum Universal Address
Trinidad and Tobago	TT	TTO	Spectrum Enterprise Geocoding (Latin America) Spectrum Universal Address
Tunisia	TN	TUN	Spectrum Enterprise Geocoding (Africa) Spectrum Universal Address

ISO Country Name	ISO 3166-1 Alpha-2	ISO 3166-1 Alpha-3	Supported Modules
Turkey	TR	TUR	Spectrum Enterprise Geocoding Spectrum Universal Address
Turkmenistan	TM	TKM	Spectrum Universal Address
Turks And Caicos Islands	TC	TCA	Spectrum Universal Address
Tuvalu	TV	TUV	Spectrum Universal Address
Uganda	UG	UGA	Spectrum Enterprise Geocoding (Africa) Spectrum Universal Address
Ukraine	UA	UKR	Spectrum Enterprise Geocoding Spectrum Universal Address
United Arab Emirates	AE	ARE	Spectrum Enterprise Geocoding (Middle East) Spectrum Universal Address
United Kingdom	GB	GBR	Spectrum Enterprise Geocoding Spectrum Spatial Routing Spectrum Universal Address
United States	US	USA	Spectrum Enterprise Geocoding Spectrum Spatial Routing Spectrum Universal Address
United States Minor Outlying Islands	UM	UMI	Spectrum Universal Address
Uruguay	UY	URY	Spectrum Enterprise Geocoding Spectrum Universal Address

ISO Country Name	ISO 3166-1 Alpha-2	ISO 3166-1 Alpha-3	Supported Modules
Uzbekistan	UZ	UZB	Spectrum Universal Address
Vanuatu	VU	VUT	Spectrum Universal Address
Venezuela, Bolivarian Republic Of	VE	VEN	Spectrum Enterprise Geocoding Spectrum Universal Address
Viet Nam	VN	VNM	Spectrum Enterprise Geocoding Spectrum Universal Address
Virgin Islands, British	VG	VGB	Spectrum Universal Address
Virgin Islands, U.S.	VI	VIR	Spectrum Universal Address
Wallis and Futuna	WF	WLF	Spectrum Universal Address
Western Sahara	EH	ESH	Spectrum Universal Address
Yemen	YE	YEM	Spectrum Enterprise Geocoding (Middle East) Spectrum Universal Address
Zambia	ZM	ZMB	Spectrum Enterprise Geocoding (Africa) Spectrum Universal Address
Zimbabwe	ZW	ZWE	Spectrum Enterprise Geocoding (Africa) Spectrum Universal Address



2 Blue Hill Plaza, #1563
Pearl River, NY 10965
USA

www.precisely.com

© 2007, 2021 Precisely. All rights reserved.