

# Spectrum™ Technology Platform

Version 2018.2.0

Data Quality Guide



# 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	64
Match Rules	67
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	112
Dataflow Templates for Matching	127

## 5 - Deduplication

---

Filtering Out Duplicate Records	136
Creating a Best of Breed Record	140

## 6 - Exception Records

---

Designing a Dataflow to Handle Exceptions	146
Designing a Dataflow for Real-Time Revalidation	147

## 7 - Lookup Tables

---

Introduction to Lookup Tables	151
Data Normalization Module Tables	151
Universal Name Module Tables	156
Viewing the Contents of a Lookup Table	158
Adding a Term to a Lookup Table	159
Removing a Term from a Lookup Table	159
Modifying the Standardized Form of a Term	159
Reverting Table Customizations	160
Creating a Lookup Table	160
Importing Data	161

## 8 - Stages Reference

---

Advanced Matching Module	165
Business Steward Module	225
Data Normalization Module	267
Universal Name Module	285

# 9 - ISO Country Codes and Module Support

---

ISO Country Codes and Module Support      320

# 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 Business Steward Module is a set of 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:

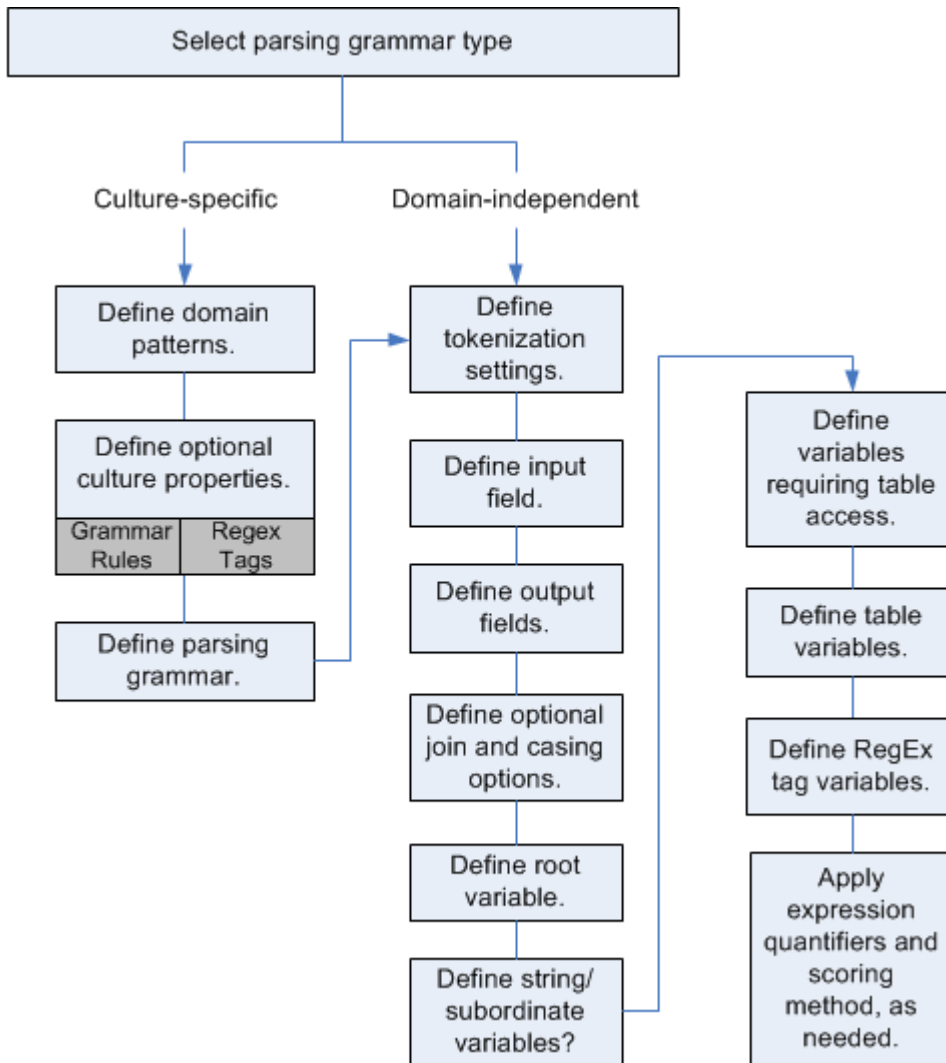
- Your input data contains names in a single field that you want to separate into given name and family name.
- Your input data contains addresses from several cultures and you want to extract address data for a specific culture only.
- Your input data includes free-form text that contains embedded email 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 274. 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
Albanian (Albania)	sq-AL
Arabic	ar

Language (Culture/Region)	Culture Code
---------------------------	--------------

---

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
-----------------------	-------

---

Arabic (Syria)	ar-SY
----------------	-------

---

Arabic (Tunisia)	ar-TN
------------------	-------

---

Language (Culture/Region)	Culture Code
---------------------------	--------------

---

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
-----------	----

---

Bulgarian (Bulgaria)	bg-BG
----------------------	-------

---

Catalan	ca
---------	----

---

Language (Culture/Region)	Culture Code
---------------------------	--------------

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
-------	----

Czech (Czech Republic)	cs-CZ
------------------------	-------

Danish	da
--------	----

Language (Culture/Region)	Culture Code
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
English (Jamaica)	en-JM
English (New Zealand)	en-NZ



Language (Culture/Region)	Culture Code
---------------------------	--------------

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
--------------	-------

Finnish	fi
---------	----

Finnish (Finland)	fi-FI
-------------------	-------

Language (Culture/Region)	Culture Code
---------------------------	--------------

---

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
--------	----

---

German (Austria)	de-AT
------------------	-------

---

German (Germany)	de-DE
------------------	-------

---

Language (Culture/Region)	Culture Code
---------------------------	--------------

---

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
-----------	----

---

Hungarian (Hungary)	hu-HU
---------------------	-------

---

Icelandic	is
-----------	----

---

Language (Culture/Region)	Culture Code
---------------------------	--------------

---

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
---------------------	-------

---

Konkani	kok
---------	-----

---

Konkani (India)	kok-IN
-----------------	--------

---

Language (Culture/Region)	Culture Code
---------------------------	--------------

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
---------------------------	-------

Malay (Malaysia)	ms-MY
------------------	-------

Marathi	mr
---------	----

Language (Culture/Region)	Culture Code
---------------------------	--------------

---

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
---------	----

---

Punjabi (India)	pa-IN
-----------------	-------

---

Romanian	ro
----------	----

---

Language (Culture/Region)	Culture Code
---------------------------	--------------

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
----------------------	-------

Spanish	es
---------	----

Spanish (Argentina)	es-AR
---------------------	-------

Language (Culture/Region)	Culture Code
---------------------------	--------------

---

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
------------------	-------

---

Spanish (Paraguay)	es-PY
--------------------	-------

---

Spanish (Peru)	es-PE
----------------	-------

---



Language (Culture/Region)	Culture Code
---------------------------	--------------

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
----------------	--------

Tamil	ta
-------	----

Tamil (India)	ta-IN
---------------	-------

Language (Culture/Region)	Culture Code
---------------------------	--------------

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
-----------------	-------

Uzbek	uz
-------	----

Uzbek (Uzbekistan, Cyrillic)	uz-Cyrl-UZ
------------------------------	------------

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

## Grammars

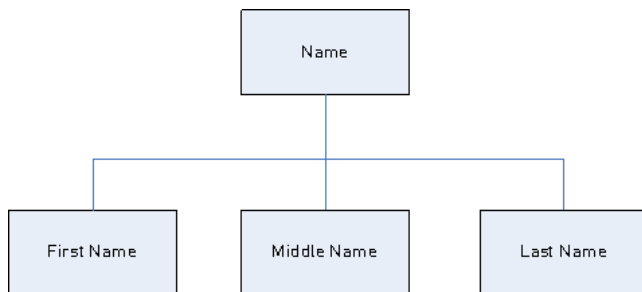
A valid parsing grammar contains:

- 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 backslash.

## 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.
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.

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

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

7. 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.
6. 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.

7. In the **Information** list, select **Step through parsing events**.

8. 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.

9. Click **Show scores** to display parser scores for root expressions, variable expressions, and the resulting matches and non-matches.

10. In the **Zoom** field, select the size of the tree view.

11. 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.

12. 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.

13. 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 multi-record name consolidation.

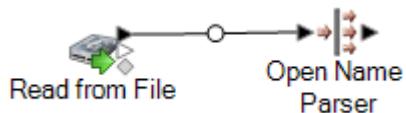
1. If you have not already done so, load the following tables onto the Spectrum™ Technology Platform server:

- Open Parser Base
- Open Parser Enhanced Names

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

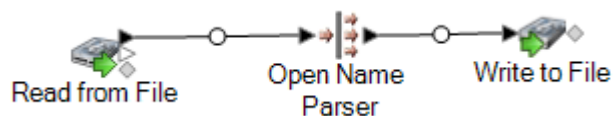
2. In Enterprise Designer, create a new dataflow.
3. Drag a source stage onto the canvas.
4. Double-click the source stage and configure it. See the *Dataflow Designer's Guide* for instructions on configuring source stages.
5. Drag an Open Name Parser stage onto the canvas and connect it to the source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

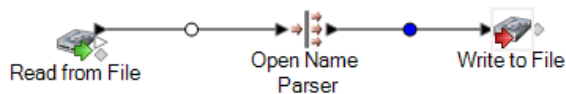
## Parsing English Names

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 Module
- 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 part of the Universal Naming Module. 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.

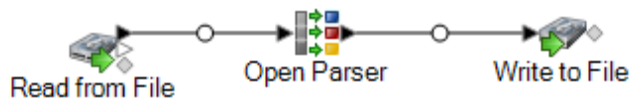
In order to understand the Arabic naming system, you search for and find these resources on the internet that explain the Arabic naming system:

- [en.wikipedia.org/wiki/Arabic\\_names](http://en.wikipedia.org/wiki/Arabic_names)
- [heraldry.sca.org/laurel/names/arabic-naming2.htm](http://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 the Data Normalization Module.

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](http://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 the Data Normalization Module.

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
刘耕宏	耕宏	刘
许惠欣	惠欣	许
蔡依林	依林	蔡
水沫	沫	水
蕭亞軒	亞軒	蕭
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

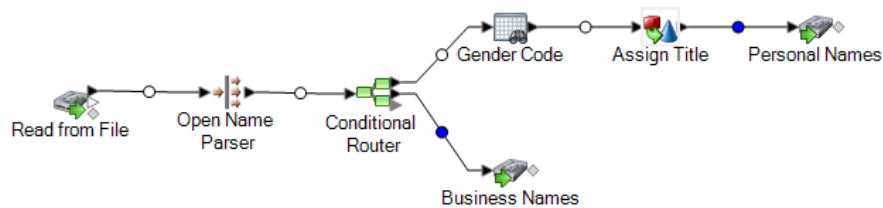
This template demonstrates how to parse mixed-culture names, such as Spanish and German names, into component parts. The parsing rule separates each token in the **Name** field and copies each token to the fields defined in the Personal and Business Names parsing grammar. For more information about this 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 the Data Normalization Module.

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 `pb.com`, "pb" 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 the Data Normalization Module.

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\Pitney Bowes\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.

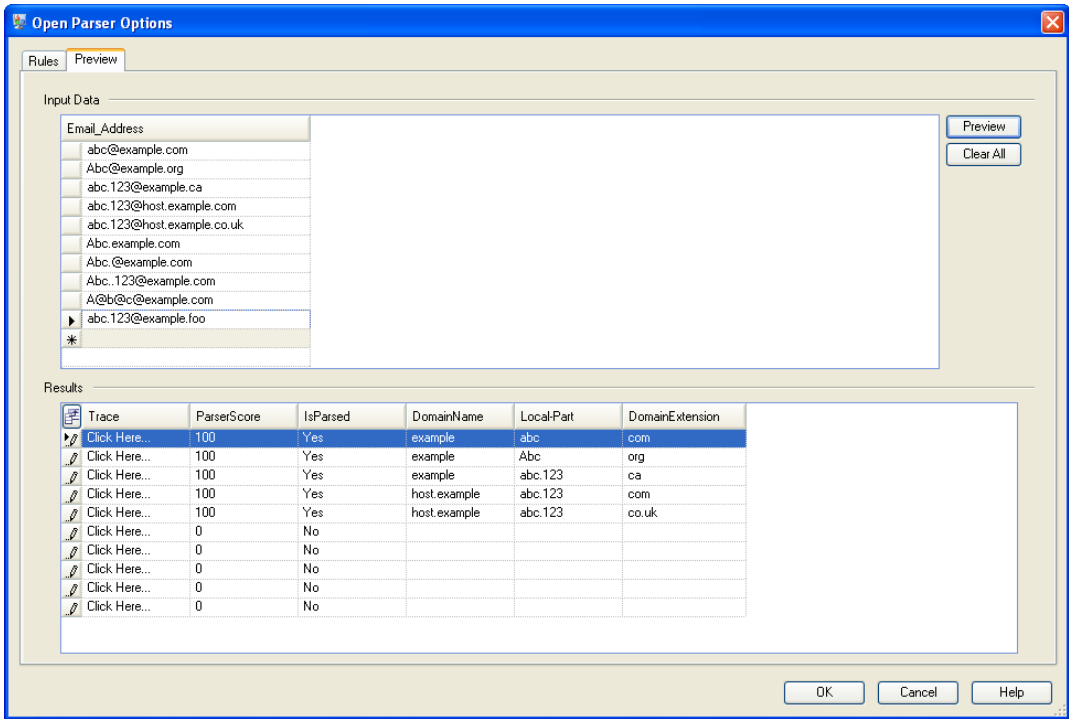
The <alphanum> variable definition is a regular expression that means any string of characters from A to Z, a to a, and 0-9. The <alphanum> variable is used throughout this parsing grammar and is defined once on the last line of the parsing grammar.

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.

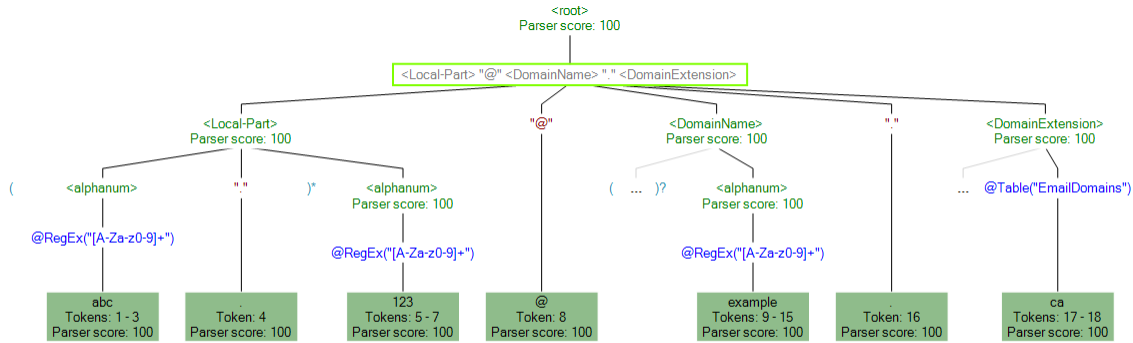
To test the parsing grammar, click the Preview tab. Type the e-mail addresses shown below in the **Email Address** field and then click **Preview**.



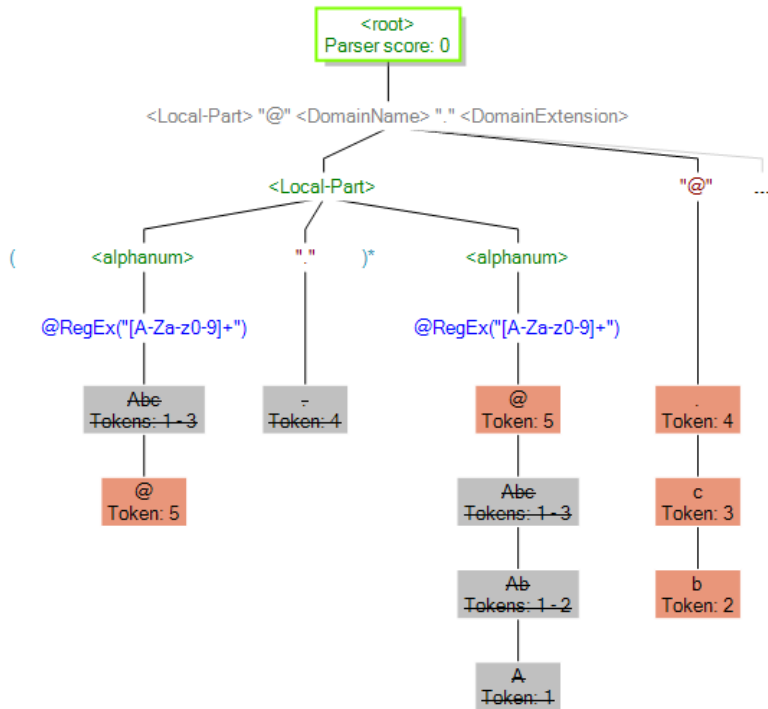
You can also type other e-mail addresses to see how the input data is parsed.

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 (Abe.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 the Data Normalization Module.

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
1404286-7534	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

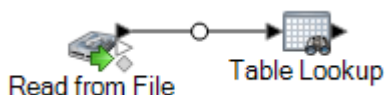
## Standardizing Terms

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 Module 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.

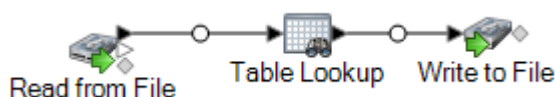
**Note:** If you do not see the table you need, contact your system administrator. The Data Normalization Module database must be loaded.

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.

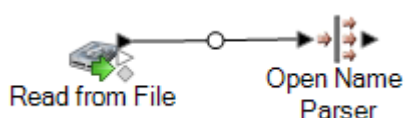
**Note:** Before beginning, make sure that your input data has a field named "Name" that contains the full name of the person.

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 Module's database load utility to load these tables. For instructions on loading tables, see the *Installation Guide*.

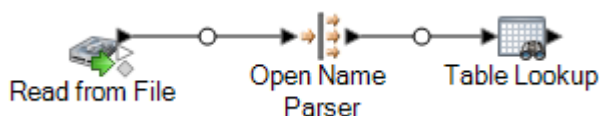
2. In Enterprise Designer, create a new dataflow.
3. Drag a source stage onto the canvas.
4. Double-click the source stage and configure it. See the *Dataflow Designer's Guide* for instructions on configuring source stages.
5. Drag an Open Name Parser stage onto the canvas and connect it to the source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

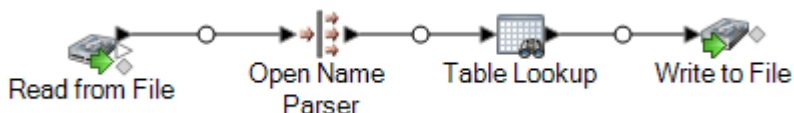


7. Double-click the Table Lookup stage on the canvas.
8. In the **Source** field, select **FirstName**.
9. In the **Destination** field, select **FirstName**.

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

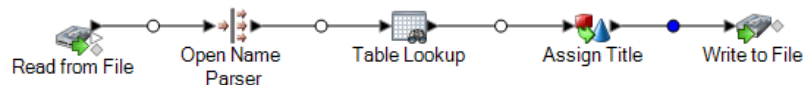
## Formalizing Personal Names

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 Data Normalization Module and the Universal Name Module.

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	64
Match Rules	67
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	112
Dataflow Templates for Matching	127

## Matching Terminology

<b>Average Score</b>	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.
<b>Baseline</b>	The selected match result that will be compared against another match result.
<b>Candidate Group</b>	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.
<b>Candidate Records</b>	All non-suspect records in a match group or candidate group.
<b>Drop</b>	A decrease in duplicates.
<b>Detail Match Record</b>	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.
<b>Duplicate Collections</b>	A duplicate collection consists of a Suspect and its Duplicate records grouped together by a CollectionNumber. Unique records always belong to CollectionNumber 0.
<b>Duplicate Records</b>	A record that matches another record within a match group. Can be a suspect or a candidate.
<b>Express Matches</b>	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.
<b>Input Records</b>	Order of the records in the matching stage before the matching sort is performed.
<b>Interflow Match</b>	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.
<b>Intraflow Match</b>	A matching stage that locates matches between similar data records within a single input stream.
<b>Lift</b>	An increase in duplicates.
<b>Match Groups</b>	(Group By) Records grouped together either by a match key or a sliding window.
<b>Match Results</b>	(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.

<b>Match Results List</b>	List of match results of a single type that MAT can analyze in the current analysis session.
<b>Match Results Type</b>	Indicates the contents of the match results. MAT uses the match results type to determine how to use the data.
<b>Matcher Stage</b>	A stage on the canvas that performs matching routines. The matcher stages are Interflow Match, Intraflow Match, and Transactional Match
<b>Missed Match</b>	A record that was previously a suspect or duplicate but is now unique.
<b>New Match</b>	A record that was previously unique but is now a suspect or duplicate.
<b>Sliding Window</b>	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.
<b>Suspect Records</b>	A driver record that is matched against candidates within a match group or a candidate group.
<b>Transactional Match</b>	A matching stage that matches suspect records against candidate records that are returned from Candidate Finder or by an external application.
<b>Unique Records</b>	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"> <li>• Primary: The record is the selected template record in a collection.</li> <li>• Secondary: The record is not the selected template record in a collection.</li> <li>• BestOfBreed: The record is the newly created best of breed record in the collection.</li> </ul>
ExpressMatchIdentified	<p>Indicates whether the match was obtained using the express match key. The possible values are Yes or No.</p>
HasDuplicates	<p>Identifies whether the record is a duplicate of another record. The options are:</p> <ul style="list-style-type: none"> <li>• Y : The record is a suspect record and has duplicates.</li> <li>• N : The record is a suspect record and has no duplicates.</li> <li>• D : The record is a candidate record and is a duplicate of the suspect record.</li> <li>• U : The record is a candidate record but is not a duplicate of the suspect record.</li> </ul>
MatchRecordType	<p>Identifies the type of match record in a collection. The possible values are:</p> <ul style="list-style-type: none"> <li>• 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.</li> <li>• Duplicate : A record that is a duplicate of the suspect record.</li> <li>• Unique : A record that has no duplicates.</li> </ul>
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"> <li>• Suspect: A suspect record is used as input to a query.</li> <li>• Candidate: A candidate record is a result returned from a query.</li> </ul>

## 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*

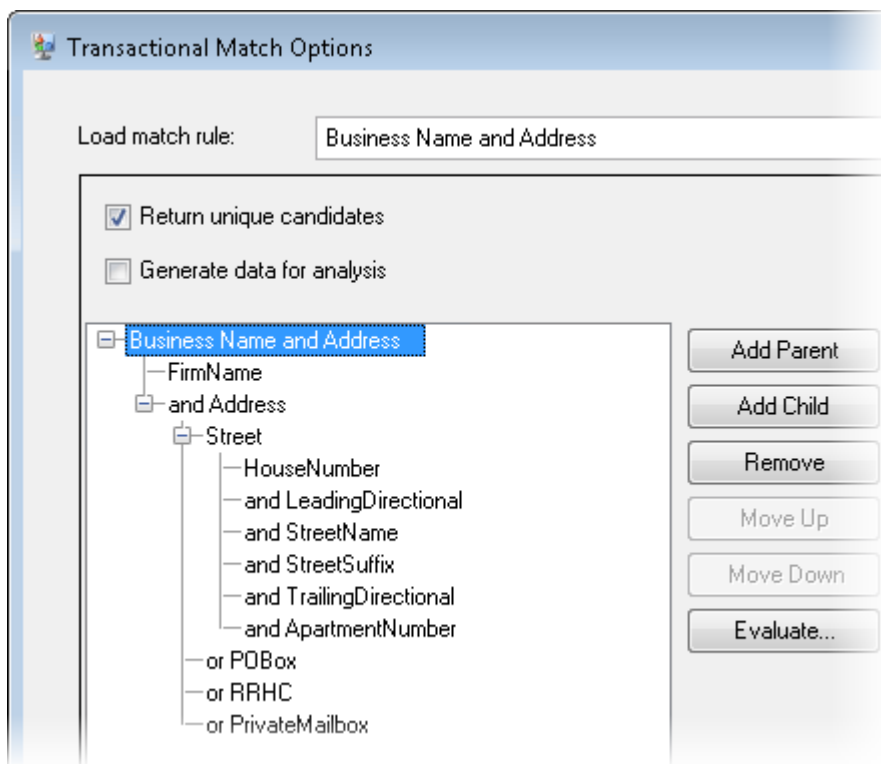
To achieve the most accurate results, you should design the match key to work well with the match rule that you will use it with. This requires you to consider how the match rule is defined.

- 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:
 

<b>All true</b>	A parent is considered a match if all children are determined to match. This method creates an "AND" connector between children.
<b>Any true</b>	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:

<b>Ignore blanks</b>	Ignores the field if it contains blank data.
<b>Count as 0</b>	Scores the field as 0 if it contains blank data.
<b>Count as 100</b>	Scores the field as 100 if it contains blank data.
<b>Compare Blanks</b>	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:

<b>Weighted Average</b>	Uses the weight of each child to determine the average match score.
<b>Average</b>	Uses the average score of each child to determine the score of a parent.
<b>Maximum</b>	Uses the highest child score to determine the score of a parent.
<b>Minimum</b>	Uses the lowest child score to determine the score of a parent.
<b>Vector Summation</b>	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:
 

<b>Ignore blanks</b>	Ignores the field if it contains blank data.
<b>Count as 0</b>	Scores the field as 0 if it contains blank data.
<b>Count as 100</b>	Scores the field as 100 if it contains blank data.
<b>Compare Blanks</b>	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** Phoenetic 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.

<b>Edit Distance</b>	Determines the similarity between two strings based on the number of deletions, insertions, or substitutions required to transform one string into another.
<b>Euclidean Distance</b>	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.
<b>Exact Match</b>	Determines if two strings are the same.
<b>Initials</b>	Used to match initials for parsed personal names.
<b>Jaro-Winkler Distance</b>	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.
<b>Keyboard Distance</b>	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).
<b>Koeln</b>	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.
<b>Kullback-Liebler Distance</b>	Determines the similarity between two strings based on the differences between the distribution of words in the two strings.
<b>Metaphone</b>	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.
<b>Metaphone (Spanish)</b>	Determines the similarity between two strings based on a phonetic representation of their characters. This option was developed to respond to limitations of Soundex.
<b>Metaphone 3</b>	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 306.

**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:

- **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:

$$\frac{(\text{numericScore} + (\text{EditDistanceScore} + \text{CharacterFrequencyScore}) / 2) / 2}$$

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.

**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

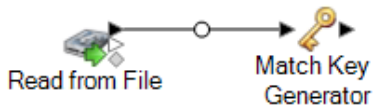
This procedure describes how to use an Intraflow Match stage to identify groups of records within a single data source (such as a file or database table) that are related to each other based on the matching criteria you specify. The dataflow groups records into collections and writes the collections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

---

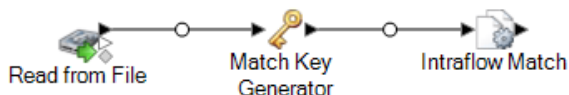
Algorithm	
-----------	--

Option Name	Description / Valid Values
	Specifies the algorithm to use to generate the match key. One of the following:
<b>Consonant</b>	Returns specified fields with consonants removed.
<b>Double Metaphone</b>	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.
<b>Koeln</b>	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.
<b>MD5</b>	A message digest algorithm that produces a 128-bit hash value. This algorithm is commonly used to check data integrity.
<b>Metaphone</b>	Returns a Metaphone coded key of selected fields. Metaphone is an algorithm for coding words using their English pronunciation.
<b>Metaphone (Spanish)</b>	Returns a Metaphone coded key of selected fields for the Spanish language. This metaphone algorithm codes words using their Spanish pronunciation.
<b>Metaphone 3</b>	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.
<b>Nysiis</b>	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.
<b>Phonix</b>	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

Option Name	Description / Valid Values
	respond to limitations of Soundex; it is more complex and therefore slower than Soundex.
<b>Soundex</b>	Returns a Soundex code of selected fields. Soundex produces a fixed-length code based on the English pronunciation of a word.
<b>Substring</b>	Returns a specified portion of the selected field.
Field name	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.
Start position	Specifies the starting position within the specified field. Not all algorithms allow you to specify a start position.
Length	Specifies the length of characters to include from the starting position. Not all algorithms allow you to specify a length.
Remove noise characters	Removes all non-numeric and non-alpha characters such as hyphens, white space, and other special characters from an input field.
Sort input	Sorts all characters in an input field or all terms in an input field in alphabetical order.
<b>Characters</b>	Sorts the characters values from an input field prior to creating a unique ID.
<b>Terms</b>	Sorts each term value from an input field prior to creating a unique ID.

8. 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.

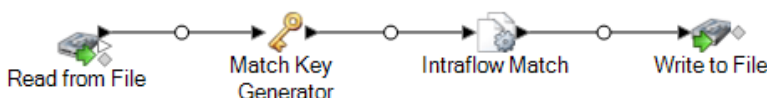
**Note:** The Dataflow Options feature in Enterprise Designer enables the match rule to be exposed for configuration at runtime.

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 68.
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 diagram shows a dataflow process for identifying household members. It begins with 'Read from File', followed by 'Match Key Generator', and 'Intraflow Match'. The output then goes to a 'Conditional Router' stage, which branches into two paths. One path goes through a 'Filter' stage, and the other path goes through a 'Stream Combiner' stage. Both paths then merge and lead to the final 'Write to File' stage.

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

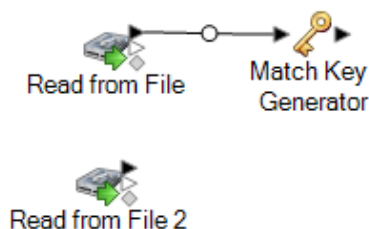
This procedure describes how to use an Interflow Match stage to identify records in one source that match records in another source. The first source contains suspect records and the second source contains candidate records. The dataflow only matches records from one source to records in another source. 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 / Valid Values
-------------	----------------------------

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

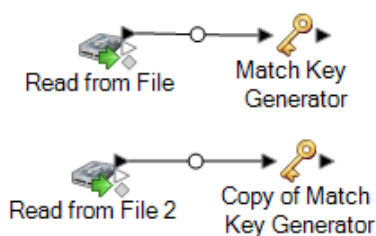


Option Name	Description / Valid Values
	Specifies the algorithm to use to generate the match key. One of the following:
<b>Consonant</b>	Returns specified fields with consonants removed.
<b>Double Metaphone</b>	Returns a code based on a phonetic representation of their characters. Double Metaphone is an improved version of the Metaphone algorithm, and attempts to account for the many irregularities found in different languages.
<b>Koeln</b>	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.
<b>MD5</b>	A message digest algorithm that produces a 128-bit hash value. This algorithm is commonly used to check data integrity.
<b>Metaphone</b>	Returns a Metaphone coded key of selected fields. Metaphone is an algorithm for coding words using their English pronunciation.
<b>Metaphone (Spanish)</b>	Returns a Metaphone coded key of selected fields for the Spanish language. This metaphone algorithm codes words using their Spanish pronunciation.
<b>Metaphone 3</b>	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.
<b>Nysiis</b>	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.
<b>Phonix</b>	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

Option Name	Description / Valid Values
	respond to limitations of Soundex; it is more complex and therefore slower than Soundex.
	<b>Soundex</b> Returns a Soundex code of selected fields. Soundex produces a fixed-length code based on the English pronunciation of a word.
	<b>Substring</b> Returns a specified portion of the selected field.
Field name	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.
Start position	Specifies the starting position within the specified field. Not all algorithms allow you to specify a start position.
Length	Specifies the length of characters to include from the starting position. Not all algorithms allow you to specify a length.
Remove noise characters	Removes all non-numeric and non-alpha characters such as hyphens, white space, and other special characters from an input field.
Sort input	Sorts all characters in an input field or all terms in an input field in alphabetical order. <b>Characters</b> Sorts the characters values from an input field prior to creating a unique ID. <b>Terms</b> Sorts each term value from an input field prior to creating a unique ID.

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.

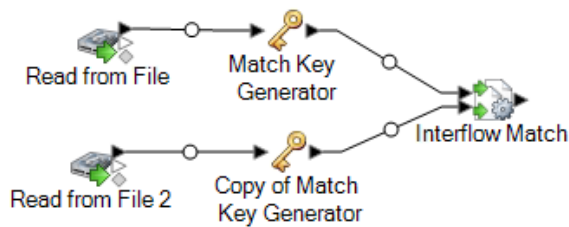
For example, if you are using Read from File input stages your dataflow would now look like this:



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.

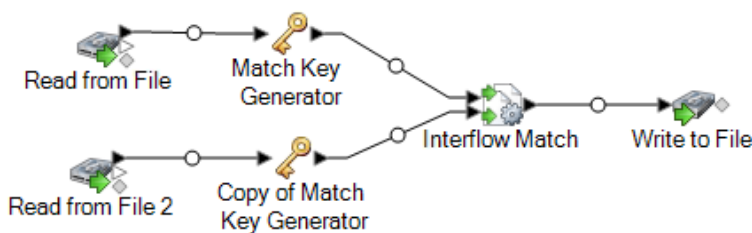
**Note:** The Dataflow Options feature in Enterprise Designer enables the match rule to be exposed for configuration at runtime.

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 68.
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.

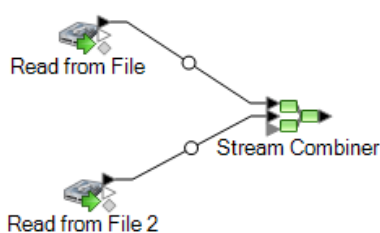
## Matching Records Between and Within Sources

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.

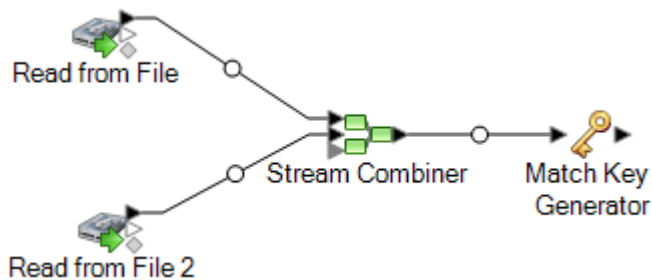
1. In Enterprise Designer, create a new dataflow.
2. Drag a source stage onto the canvas.
3. Double-click the source stage and configure it. See the *Dataflow Designer's Guide* for instructions on configuring source stages.
4. Drag a second source stage onto the canvas and configure it to read the second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

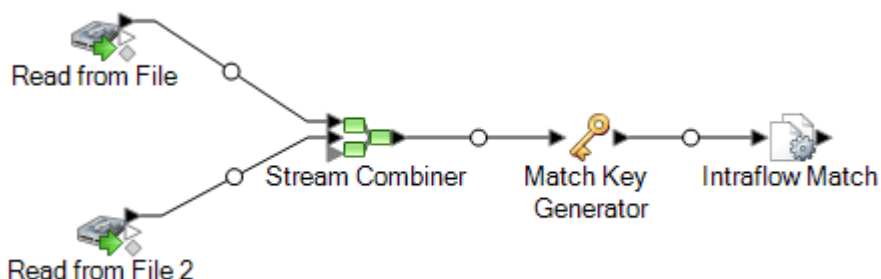
Option Name	Description / Valid Values
Algorithm	

Option Name	Description / Valid Values
	Specifies the algorithm to use to generate the match key. One of the following:
<b>Consonant</b>	Returns specified fields with consonants removed.
<b>Double Metaphone</b>	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.
<b>Koeln</b>	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.
<b>MD5</b>	A message digest algorithm that produces a 128-bit hash value. This algorithm is commonly used to check data integrity.
<b>Metaphone</b>	Returns a Metaphone coded key of selected fields. Metaphone is an algorithm for coding words using their English pronunciation.
<b>Metaphone (Spanish)</b>	Returns a Metaphone coded key of selected fields for the Spanish language. This metaphone algorithm codes words using their Spanish pronunciation.
<b>Metaphone 3</b>	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.
<b>Nysiis</b>	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.
<b>Phonix</b>	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

Option Name	Description / Valid Values
	respond to limitations of Soundex; it is more complex and therefore slower than Soundex.
	<b>Soundex</b> Returns a Soundex code of selected fields. Soundex produces a fixed-length code based on the English pronunciation of a word.
	<b>Substring</b> Returns a specified portion of the selected field.
Field name	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.
Start position	Specifies the starting position within the specified field. Not all algorithms allow you to specify a start position.
Length	Specifies the length of characters to include from the starting position. Not all algorithms allow you to specify a length.
Remove noise characters	Removes all non-numeric and non-alpha characters such as hyphens, white space, and other special characters from an input field.
Sort input	Sorts all characters in an input field or all terms in an input field in alphabetical order. <b>Characters</b> Sorts the characters values from an input field prior to creating a unique ID. <b>Terms</b> Sorts each term value from an input field prior to creating a unique ID.

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.

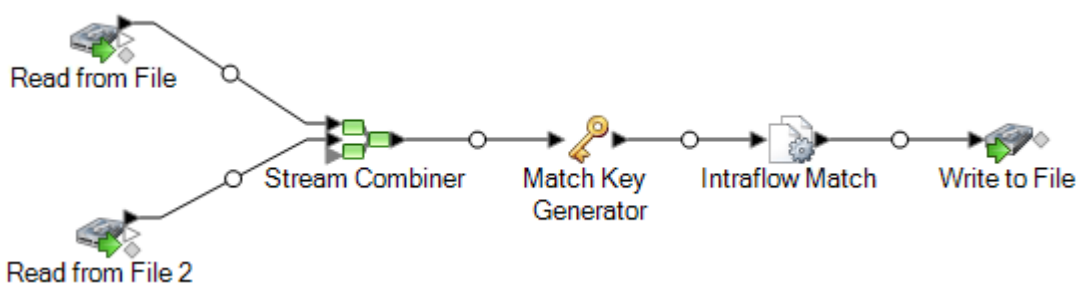
**Note:** The Dataflow Options feature in Enterprise Designer enables the match rule to be exposed for configuration at runtime.

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 68.
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*.

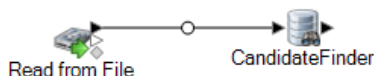
## Matching Records Against a Database

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:

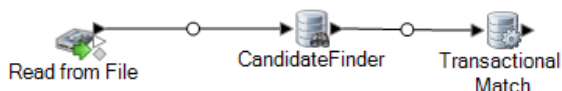
```
SELECT FirstName, LastName, Address, City, State, PostalCode
FROM Customer_Table
WHERE LastName = ${Customer_LastName};
```

8. 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.

9. Click **OK**.
10. 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:



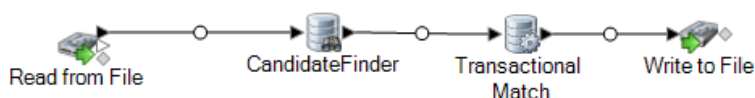
Transactional Match matches suspect records against candidate records that are returned from the Candidate Finder stage. Transactional Match uses matching rules to compare the suspect record to all candidate records with the same candidate group number (assigned in Candidate Finder) to identify duplicates.

11. Double-click the Transactional Match stage on the canvas.
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. For information about modifying the other options, see [Building a Match Rule](#) on page 68.
14. When you are done configuring the Transactional Match stage, click **OK**.
15. 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:



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

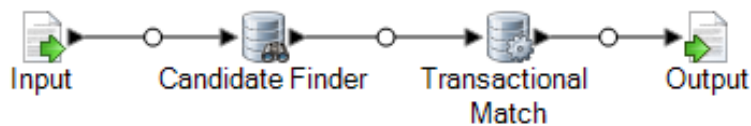
For information on configuring sink stages, see the *Dataflow Designer's Guide*.

You now have a dataflow that 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

### Download the sample dataflow

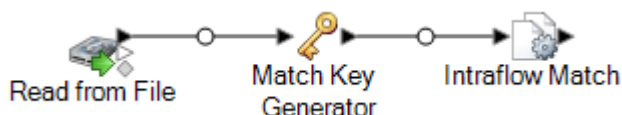
If you have records that you want to match and you want to use more than one matching operation, you can create a dataflow that uses more than one match key then combines the results to effectively match on multiple separate criteria. For example, say you want to create a dataflow that matches records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

The name and address match  
OR  
The date of birth and government ID match

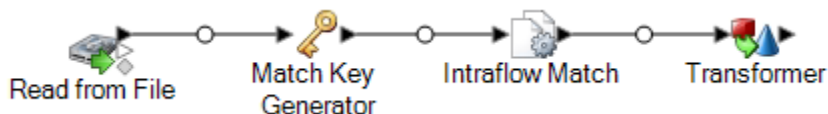
To perform matching using this logic, you create a dataflow that performs name and address matching in one stage, and date of birth and government ID matching in another stage, then combine the matching records into a single collection.

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.

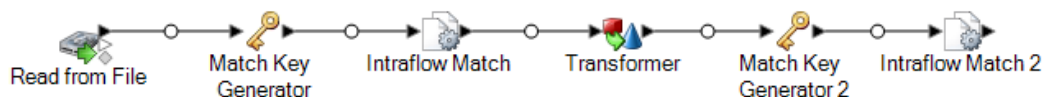
1. In Enterprise Designer, create a new dataflow.
2. Drag a source stage onto the canvas.
3. Double-click the source stage and configure it. See the *Dataflow Designer's Guide* for instructions on configuring source stages.
4. Define the first matching pass. The results of this first matching pass will 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.

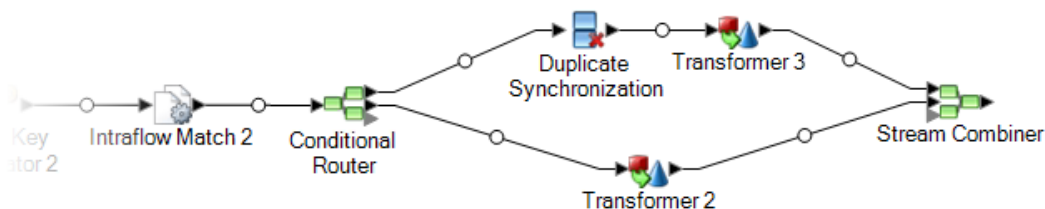
For example, if you want the second matching pass to match date of birth and government ID, you might create a match key based on the fields containing the birthday and government ID.

- 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.

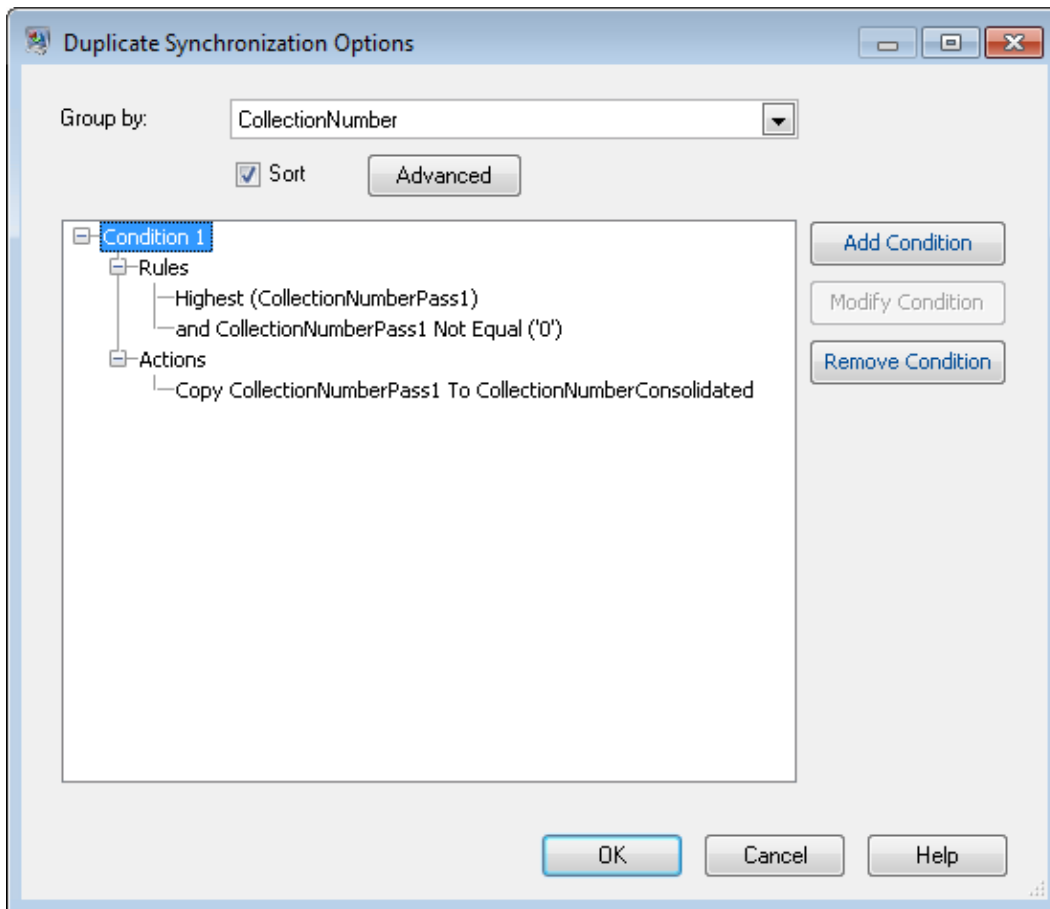
- a) Create the dataflow snippet shown below following the second Intraflow Match stage:



- 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.

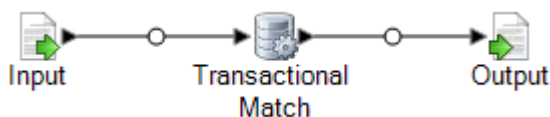
8. After the Stream Combiner you will have collections of records that match in either of the matching passes. The CollectionNumberConsolidated field indicates the matching records. You can add a sink or any additional processing you wish to perform after the Stream Combiner stage.

## 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.pb.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_fields>
        <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.pb.com/"

xmlns:ns3="http://www.pb.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_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 / Valid Values
Algorithm	

Option Name	Description / Valid Values
	Specifies the algorithm to use to generate the match key. One of the following:
<b>Consonant</b>	Returns specified fields with consonants removed.
<b>Double Metaphone</b>	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.
<b>Koeln</b>	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.
<b>MD5</b>	A message digest algorithm that produces a 128-bit hash value. This algorithm is commonly used to check data integrity.
<b>Metaphone</b>	Returns a Metaphone coded key of selected fields. Metaphone is an algorithm for coding words using their English pronunciation.
<b>Metaphone (Spanish)</b>	Returns a Metaphone coded key of selected fields for the Spanish language. This metaphone algorithm codes words using their Spanish pronunciation.
<b>Metaphone 3</b>	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.
<b>Nysiis</b>	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.
<b>Phonix</b>	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

Option Name	Description / Valid Values
	respond to limitations of Soundex; it is more complex and therefore slower than Soundex.
<b>Soundex</b>	Returns a Soundex code of selected fields. Soundex produces a fixed-length code based on the English pronunciation of a word.
<b>Substring</b>	Returns a specified portion of the selected field.
Field name	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.
Start position	Specifies the starting position within the specified field. Not all algorithms allow you to specify a start position.
Length	Specifies the length of characters to include from the starting position. Not all algorithms allow you to specify a length.
Remove noise characters	Removes all non-numeric and non-alpha characters such as hyphens, white space, and other special characters from an input field.
Sort input	Sorts all characters in an input field or all terms in an input field in alphabetical order. <ul style="list-style-type: none"> <li><b>Characters</b>      Sorts the characters values from an input field prior to creating a unique ID.</li> <li><b>Terms</b>            Sorts each term value from an input field prior to creating a unique ID.</li> </ul>

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 information about the matching processes in a dataflow, such as the number of duplicate records, the average match score, and so on. 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:

<b>Input Records</b>	The total number of records processed by the matcher stage.
<b>Unique Records</b>	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.
<b>Match Groups</b>	(Group By) Records grouped together either by a match key or a sliding window.
<b>Duplicate Collections</b>	A duplicate collection consists of a Suspect and its Duplicate records grouped together by a CollectionNumber. Unique records always belong to CollectionNumber 0.
<b>Express Matches</b>	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.
<b>Average Score</b>	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:

<b>Duplicate Collections</b>	A duplicate collection consists of a Suspect and its Duplicate records grouped together by a CollectionNumber. Unique records always belong to CollectionNumber 0.
<b>Express Matches</b>	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.

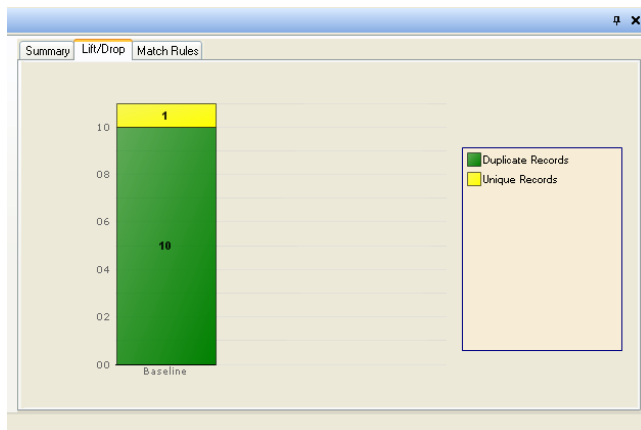
<b>Average Score</b>	The average match score of all duplicates. The possible values are 0-100, with 0 indicating a poor match and 100 indicating an exact match.
<b>Input Suspects</b>	The number of records in the input stream that the matcher tried to match to other records.
<b>Suspects with Duplicates</b>	The number of input suspects that matched at least one candidate record.
<b>Unique Suspects</b>	The number of input suspects that did not match any candidate records.
<b>Suspects with Candidates</b>	The number of input suspects that had at least one candidate record in its match group and therefore had at least one match attempt.
<b>Suspects without Candidates</b>	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:

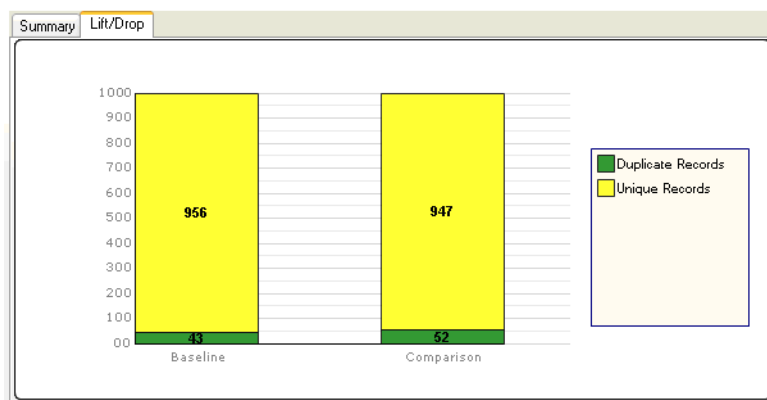
<b>Average Score</b>	The average match score of all duplicates. The possible values are 0-100, with 0 indicating a poor match and 100 indicating an exact match.
<b>Input Suspects</b>	The number of records in the input stream that the matcher tried to match to other records.
<b>Suspects with Duplicates</b>	The number of input suspects that matched at least one candidate record.
<b>Unique Suspects</b>	The number of input suspects that did not match any candidate records.
<b>Suspects with Candidates</b>	The number of input suspects that had at least one candidate record in its match group and therefore had at least one match attempt.
<b>Suspects without Candidates</b>	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:



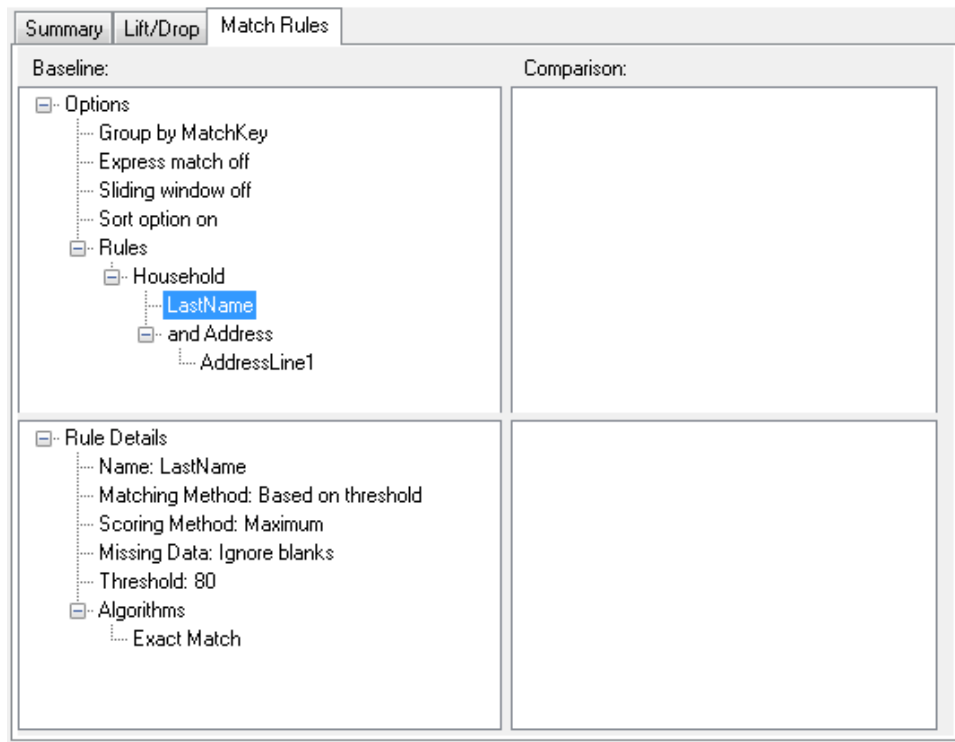
If both a baseline and a comparison job are selected, a chart for the baseline and comparison jobs are shown side by side:



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.

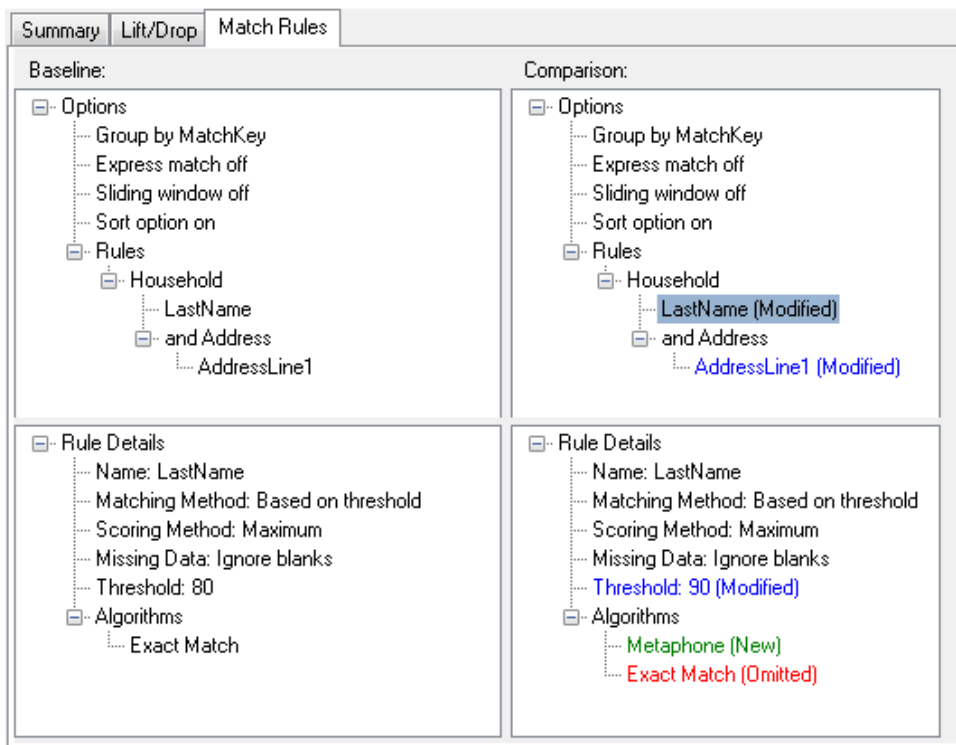
To view rule details, select a node in the hierarchy.



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.

For example:



## 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 <b>Remove</b> .		X	X

For information about the match rate chart, see [Match Rate Chart](#) on page 120.

- 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:

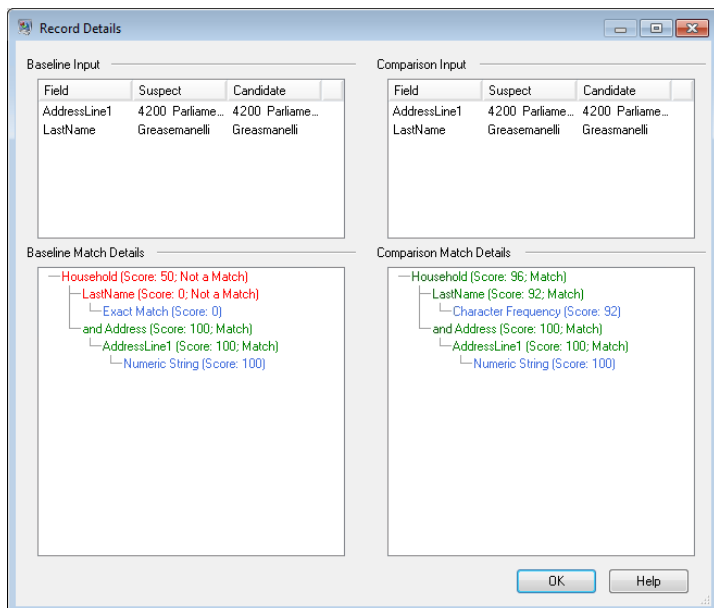
- New Matches: (Intraflow) Displays all new matches and its related suspects. This view combines the results of Suspects with New Duplicates and New Suspects into one view.
  - New Matched Suspects: (Interflow and Transactional Match) Displays suspects that had no duplicates in the baseline but have 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** window 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.

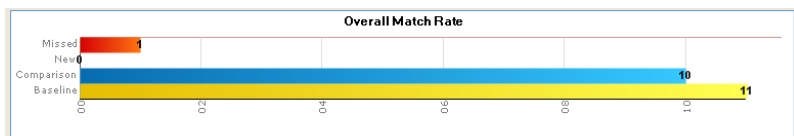
- Baseline Input—Displays the field level data, from both the suspect and candidate, used in the match.
- Baseline Match Details—Displays scoring information for each node in the match rules.
- Comparison Input—Displays the field level data, from both the suspect and candidate, used in the match.

- 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:

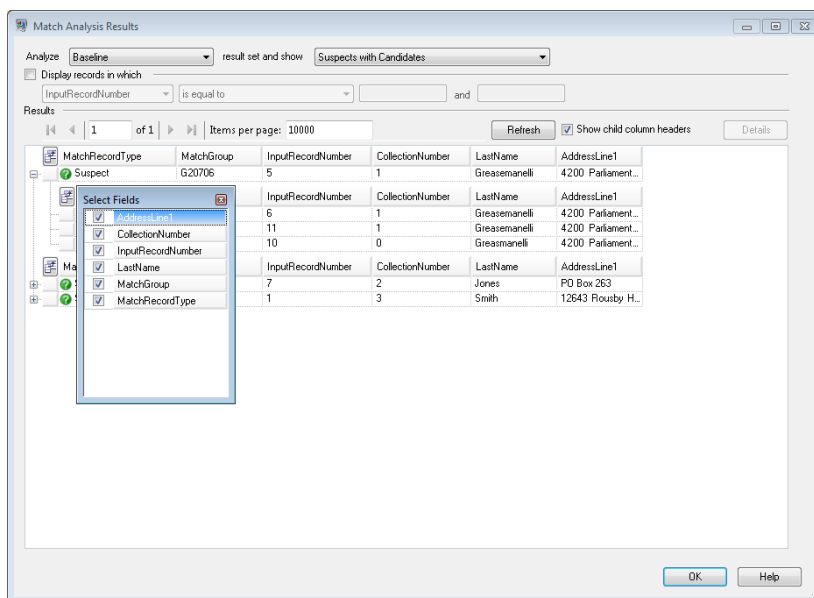
- Overall Match Rate
- Baseline Matches: Total number of matches in the baseline result.
- Comparison Matches: Total number of matches in the comparison result.
- New Matches: A count of all records that were unique in the baseline result, but are a suspect or duplicate in the comparison result.



- 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.

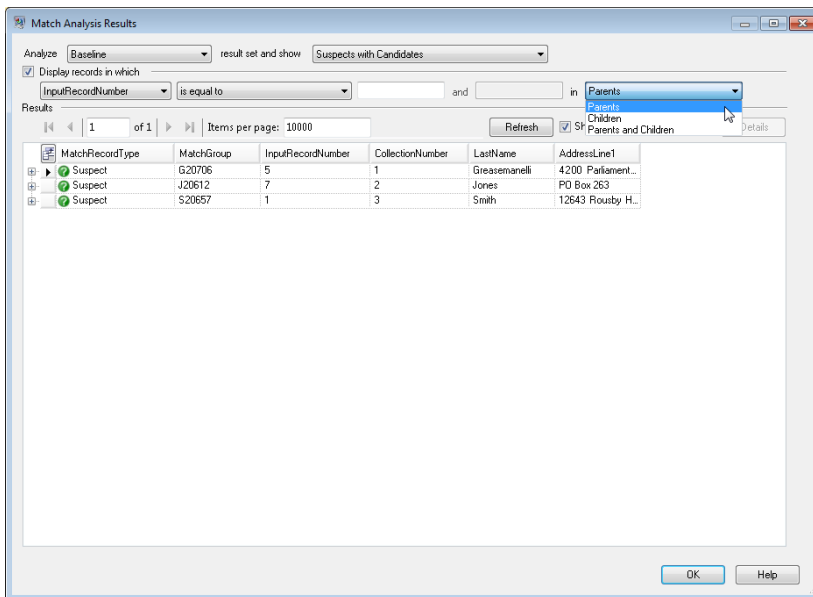
The operators you can choose are:

- 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's 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 dataflow's 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.

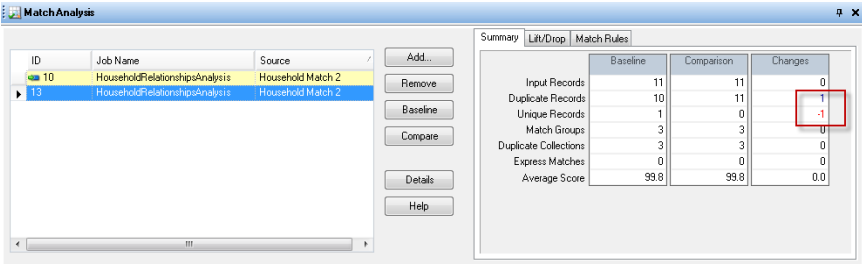
**Tip:** If there are a large number of dataflows and you want to filter the dataflows, select a filter option from the **Show only jobs where** drop-down list.

- 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

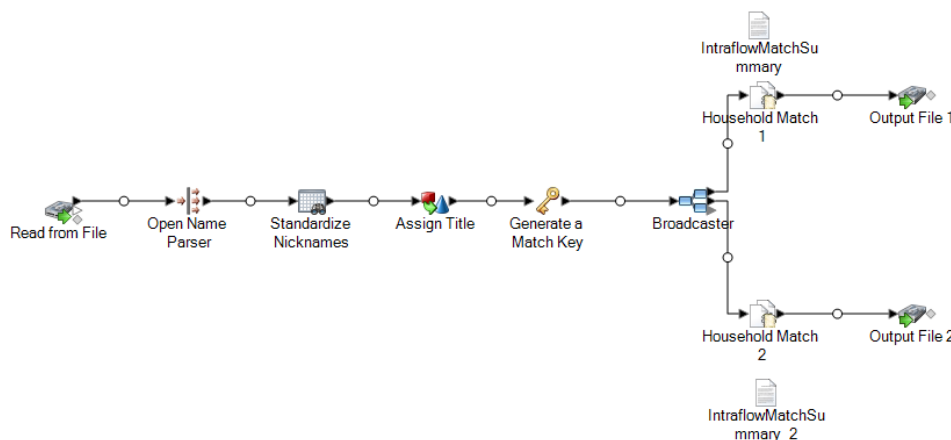
To remove a match results from the Match Results List, select a match results in the Match Results List and then click **Remove**.

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

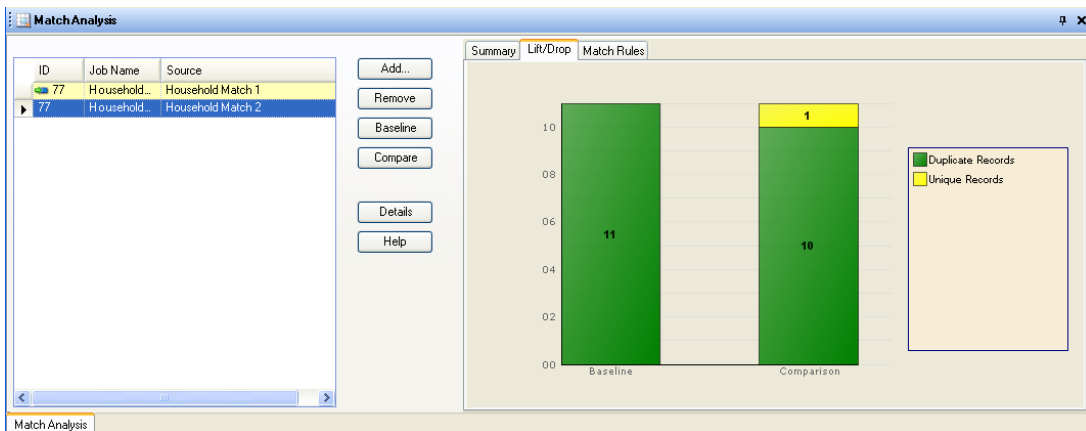
This example demonstrates how to use the Match Analysis tool to compare the lift/drop rates of two different matches. Before the data is sent through a matcher, it is split into two streams using a Broadcaster. Each stream is then sent through an Intraflow Match stage. Each data stream includes identical 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 modules: Advanced Matching Module, Data Normalization Module, and Universal Name Module. It also requires you to load the Table Lookup core database and the Open Parser base tables.

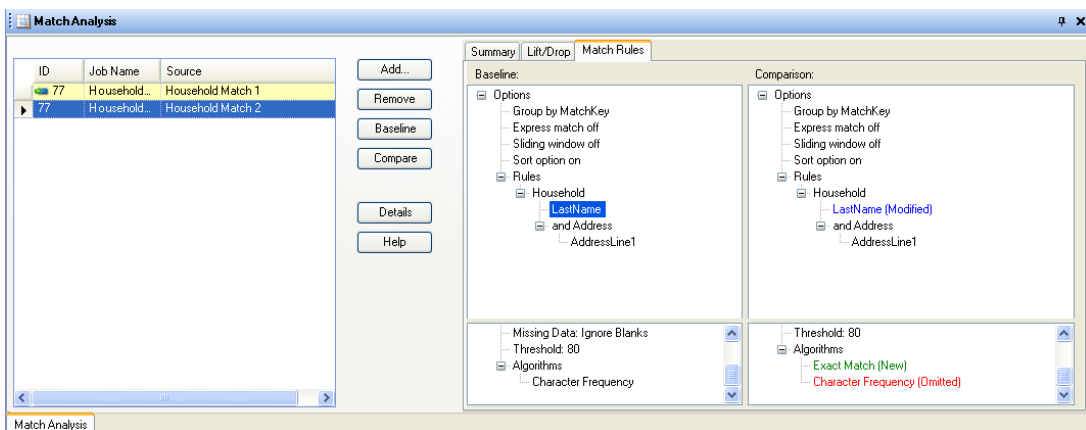
To use view this example:

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**. The Summary Results display.
5. Click the **Lift/Drop** tab. The Lift/Drop chart displays.



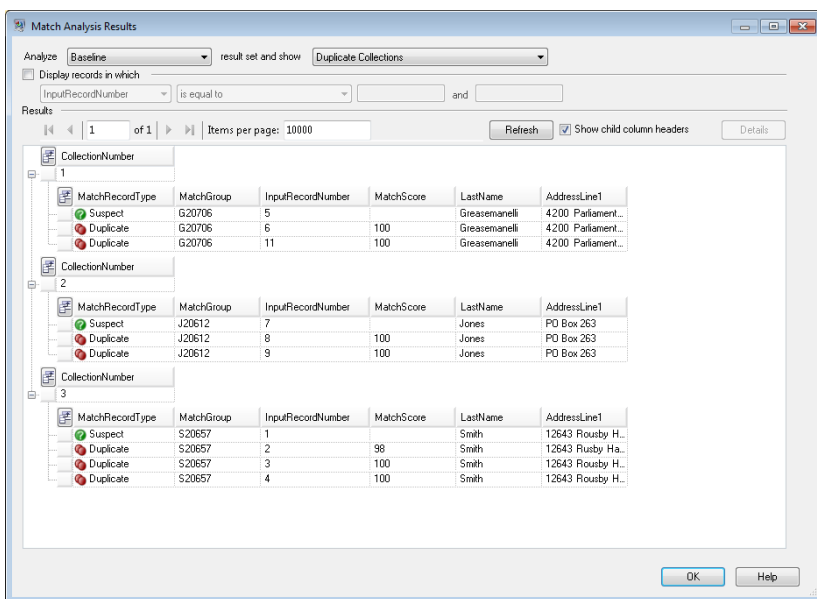
This chart shows the differences between the duplicate and unique records generated for the different match rules used.

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



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

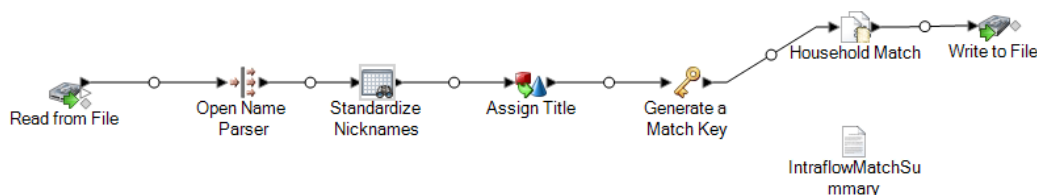
### Identifying Members of a Household

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 modules: Advanced Matching Module, Data Normalization Module, and Universal Name Module

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.

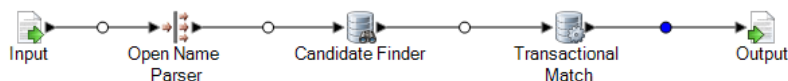
## Determining if a Prospect is a Customer

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 the Advanced Matching Module and Universal Name Module.

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	136
Creating a Best of Breed Record	140

## 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

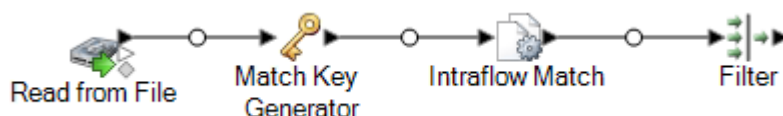
[Matching Records from One Source to Another Source](#) on page 86

[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: <b>Non-Numeric</b> Choose this option if the field contains non-numeric data (for example, string data). <b>Numeric</b> 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:
<b>Contains</b>	Determines if the field contains the value specified. For example, "sailboat" contains the value "boat".
<b>Equal</b>	Determines if the field contains the exact value specified.
<b>Greater Than</b>	Determines if the field value is greater than the value specified. This operation only works on numeric fields.
<b>Greater Than Or Equal To</b>	Determines if the field value is greater than or equal to the value specified. This operation only works on numeric fields.
<b>Highest</b>	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.
<b>Is Empty</b>	Determines if the field contains no value.
<b>Is Not Empty</b>	Determines if the field contains any value.
<b>Less Than</b>	Determines if the field value is less than the value specified. This operation only works on numeric fields.
<b>Less Than Or Equal To</b>	Determines if the field value is less than or equal to the value specified. This operation only works on numeric fields.
<b>Longest</b>	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.
<b>Lowest</b>	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.
<b>Most Common</b>	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.
<b>Not Equal</b>	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><b>Note:</b> This option is not available if you select the operator Highest, Lowest, or Longest.</p> <p><b>Field</b> Choose this option if you want to compare another dataflow field's value to the field.</p> <p><b>String</b> 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 <b>Field type</b> field, select a dataflow field. If you selected String in the <b>Value type</b> field, type the value you want to use in the comparison.</p> <p><b>Note:</b> 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

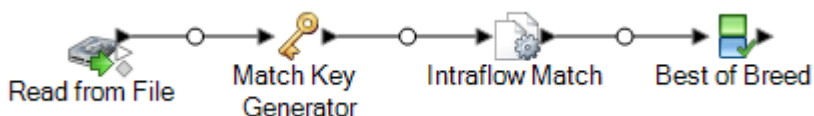
[Matching Records from One Source to Another Source](#) on page 86

[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: <b>Non-Numeric</b> Choose this option if the field contains non-numeric data (for example, string data). <b>Numeric</b> 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:
<b>Contains</b>	Determines if the field contains the value specified. For example, "sailboat" contains the value "boat".
<b>Equal</b>	Determines if the field contains the exact value specified.
<b>Greater Than</b>	Determines if the field value is greater than the value specified. This operation only works on numeric fields.
<b>Greater Than Or Equal To</b>	Determines if the field value is greater than or equal to the value specified. This operation only works on numeric fields.
<b>Highest</b>	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.
<b>Is Empty</b>	Determines if the field contains no value.
<b>Is Not Empty</b>	Determines if the field contains any value.
<b>Less Than</b>	Determines if the field value is less than the value specified. This operation only works on numeric fields.
<b>Less Than Or Equal To</b>	Determines if the field value is less than or equal to the value specified. This operation only works on numeric fields.
<b>Longest</b>	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.
<b>Lowest</b>	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.
<b>Most Common</b>	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.
<b>Not Equal</b>	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><b>Note:</b> This option is not available if you select the operator Highest, Lowest, or Longest.</p> <p><b>Field</b> Choose this option if you want to compare another dataflow field's value to the field.</p> <p><b>String</b> 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 <b>Field type</b> field, select a dataflow field. If you selected String in the <b>Value type</b> field, type the value you want to use in the comparison.</p> <p><b>Note:</b> 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><b>Field</b> Choose this option if you want to copy a value from a field to the best of breed record.</p> <p><b>String</b> 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 <b>Source data</b> field.</p>

Option	Description
Accumulate source data	<p>If the data in the <b>Source data</b> 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>

12. Click **OK**.

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

13. Click **OK** to close the **Best of Breed Options** window.

14. 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:



15. 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	146
Designing a Dataflow for Real-Time Revalidation	147

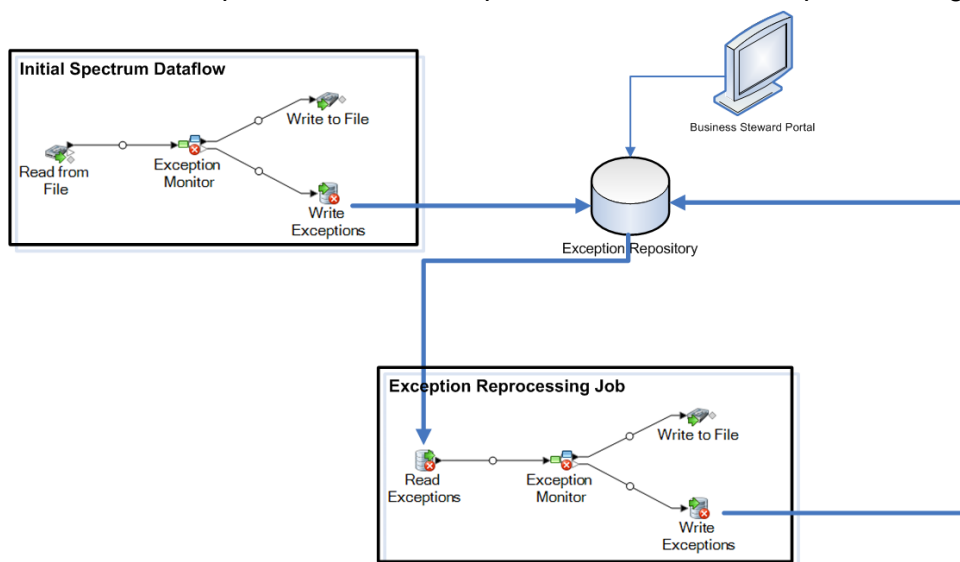
## Designing a Dataflow to Handle Exceptions

If you have licensed the Business Steward Module, 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 Business Steward Portal, a browser-based tool, which allows you to review and edit exception records. Once edited, the records are marked as "Approved", which makes the records available to be 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 Business Steward 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 226.

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 Business Steward Portal, a browser-based tool for viewing and modifying exception records. Using our example, the data steward could use the Exception Editor in the Business Steward Portal to manually add postal codes to the exception records and mark them as "Approved".

Once a record is marked as "Approved" in the Business Steward Portal, the record is available to be read back into a Spectrum™ Technology Platform dataflow. This is accomplished by using a Read Exceptions stage. If any records still result in an exception they are once 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

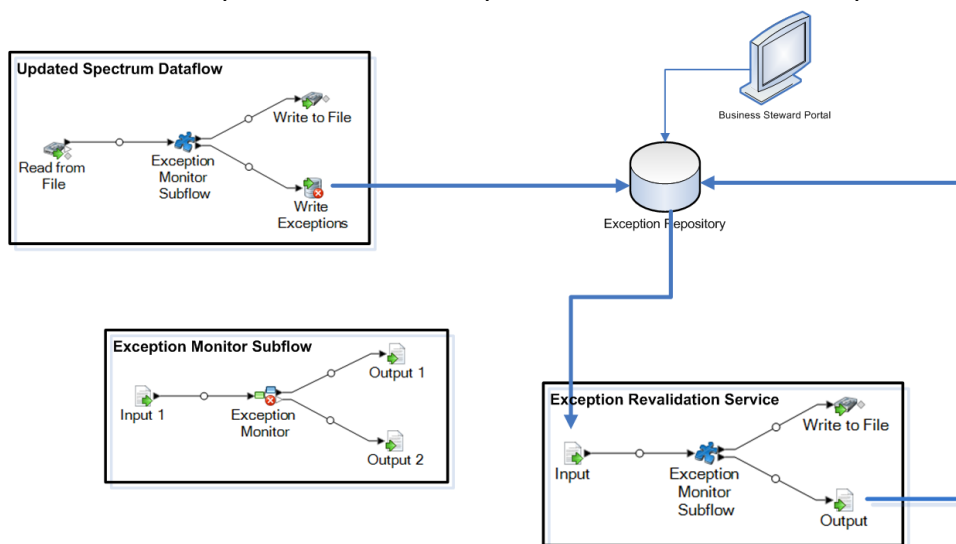
If you are using exception management in your dataflow, you can use the revalidation feature to rerun exception records through the validation process after they have been corrected in the Business Steward Portal. This enables you to determine if the change you made causes the record to process

successfully in a real-time manner; you don't 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 Business Steward 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. Continuing with our example, 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 Business Steward 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 253 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 exception records in the Business Steward Portal Exception Editor and click **Save**. Like the job, the service contains the exception monitor subflow that uses the same business logic to reprocess the record(s). 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 in the "Action after revalidation" field:

- **Reprocess records**—Records will be deleted from the repository and reprocessed.
- **Approve records**—Records will be marked as approved and sent back to the repository.

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, 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 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 Exception 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	151
Data Normalization Module Tables	151
Universal Name Module Tables	156
Viewing the Contents of a Lookup Table	158
Adding a Term to a Lookup Table	159
Removing a Term from a Lookup Table	159
Modifying the Standardized Form of a Term	159
Reverting Table Customizations	160
Creating a Lookup Table	160
Importing Data	161

## 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 Module 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 151.

- 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 151.

### *Base Tables*

Base tables are provided with the Data Normalization Module 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 Module 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 Module 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 Module 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 Module 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 Module 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 Module 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 Module installation package and thus require an additional license. For more information, contact your account executive.

Asian Plus Pack tables must be loaded using the Data Normalization Module 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 151.

### *Base Tables*

Base tables are provided with the Data Normalization Module 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 Module 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 Module 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 Module 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 Module 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 151.

### *Base Tables*

Base tables are provided with the Universal Name Module 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 Module 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 Module 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 Module 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 Module 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
<b>Find a specific term</b>	In the <b>Starts with</b> field, type the term you want to find then click <b>Refresh</b> .
<b>Page through the table</b>	Click the forward and back icons to the right of the <b>Refresh</b> button.
<b>Change the number of terms displayed per page</b>	Change the value in the <b>Items per page</b> field.
<b>View all the lookup terms for each standardized term in a Table Lookup table</b>	In the <b>View by</b> field select <b>Standardized Term (Grouping)</b> . 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 PB to Pitney Bowes, you would enter PB as the lookup term, and Pitney Bowes 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 PB and PB Software, and the standardized term is Pitney Bowes, and you want to change the standardized form to Pitney Bowes 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

The Advanced Matching Module, Data Normalization Module, and Universal Name Module 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 159.

## 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 - Stages Reference

## In this section

---

Advanced Matching Module	165
Business Steward Module	225
Data Normalization Module	267
Universal Name Module	285

# Advanced Matching Module

## Advanced Matching Module

The Advanced Matching Module matches records between and/or within any number of input files. You can also use the Advanced Matching Module to match on a variety of fields including name, address, name and address, or non-name/address fields, such as social security number or date of birth.

### Components

The Advanced Matching Module consists of:

- **Best-of-Breed**—This stage selects a best-of-breed record from the duplicates cluster by selecting a template record then using this record to build a best composite record. This record is then returned as the survivor record.
- **Candidate Finder**—This stage obtains the candidate records that will form the set of potential matches that the Transactional Match stage will evaluate. Candidate Finder is used in conjunction with Transactional Match.
- **Duplicate Synchronization**—This stage allows you to specify which fields from a collection of records will be copied (posted) to the corresponding fields of all records in the collection.
- **Filter**—This stage allows you to specify the criteria that records must satisfy in order to be retained or removed from a collection of duplicate records, either for further processing downstream or to be written to 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 give you a better understanding of the outcome of the matching process and assist you in improving your matching results.
- **Match Key Generator**—This stage creates a non-unique key shared by all like records. Records that share the same key can then be grouped together for comparison.
- **Private Match**—This stage enables two entities to compare datasets and identify common records, or customers, without compromising sensitive information. By using an encrypt function security is retained while a match function is performed, and then a decrypt 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 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 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 <b>Group by</b> 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="553 373 1427 520">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 <b>Override sort performance options</b> box then specify the values you want in these fields:</p> <p data-bbox="553 537 1427 905"><b>In memory record limit</b> 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="818 779 1427 905"><b>Note:</b> Be careful in environments where there are jobs running concurrently because increasing the <b>In memory record limit</b> setting increases the likelihood of running out of memory.</p> <p data-bbox="553 936 1427 1167"><b>Maximum number of temporary files</b> 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 data-bbox="735 1188 1427 1310" style="background-color: #f0f0f0; padding: 10px; margin: 10px 0;"> <math display="block">\frac{(NumberOfRecords \times 2)}{InMemoryRecordLimit} = NumberOfTempFiles</math> </div> <p data-bbox="735 1325 1427 1381">Note that the maximum number of temporary files cannot be more than 1,000.</p> <p data-bbox="553 1398 1427 1566"><b>Enable compression</b> Specifies that temporary files are compressed when they are written to disk.</p> <p data-bbox="639 1482 1427 1566"><b>Note:</b> The optimal sort performance settings depends on your server's hardware configuration. Nevertheless, the following equation generally produces good sort performance:</p> <div data-bbox="639 1587 1427 1688" style="background-color: #f0f0f0; padding: 10px; margin: 10px 0;"> <math display="block">(InMemoryRecordLimit \times MaxNumberOfTempFiles \div 2) \geq TotalNumberOfRecords</math> </div>
Keep original records	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.

Option Name	Description / Valid Values
Use first record	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.
Define template record	Select this option to define rules for selecting the template record. For more information, see <a href="#">Defining Template Record Rules</a> on page 168.

### 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:
	<b>Non-Numeric</b> Choose this option if the field contains non-numeric data (for example, string data).
	<b>Numeric</b> 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:
<b>Contains</b>	Determines if the field contains the value specified. For example, "sailboat" contains the value "boat".
<b>Equal</b>	Determines if the field contains the exact value specified.
<b>Greater Than</b>	Determines if the field value is greater than the value specified. This operation only works on numeric fields.
<b>Greater Than Or Equal To</b>	Determines if the field value is greater than or equal to the value specified. This operation only works on numeric fields.
<b>Highest</b>	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.
<b>Is Empty</b>	Determines if the field contains no value.
<b>Is Not Empty</b>	Determines if the field contains any value.
<b>Less Than</b>	Determines if the field value is less than the value specified. This operation only works on numeric fields.
<b>Less Than Or Equal To</b>	Determines if the field value is less than or equal to the value specified. This operation only works on numeric fields.
<b>Longest</b>	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.
<b>Lowest</b>	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.
<b>Most Common</b>	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.
<b>Not Equal</b>	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><b>Note:</b> This option is not available if you select the operator Highest, Lowest, or Longest.</p> <p><b>Field</b> Choose this option if you want to compare another dataflow field's value to the field.</p> <p><b>String</b> 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 <b>Field type</b> field, select a dataflow field. If you selected String in the <b>Value type</b> field, type the value you want to use in the comparison.</p> <p><b>Note:</b> This option is not available if you select the operator Highest, Lowest, or Longest.</p>

5. Click **OK**.
6. 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

1. In the Best of Breed stage, under **Best of Breed Settings**, click the **Rules** node in the tree.
2. Click **Add Rule**.
3. 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: <b>Non-Numeric</b> Choose this option if the field contains non-numeric data (for example, string data). <b>Numeric</b> 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:
<b>Contains</b>	Determines if the field contains the value specified. For example, "sailboat" contains the value "boat".
<b>Equal</b>	Determines if the field contains the exact value specified.
<b>Greater Than</b>	Determines if the field value is greater than the value specified. This operation only works on numeric fields.
<b>Greater Than Or Equal To</b>	Determines if the field value is greater than or equal to the value specified. This operation only works on numeric fields.
<b>Highest</b>	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.
<b>Is Empty</b>	Determines if the field contains no value.
<b>Is Not Empty</b>	Determines if the field contains any value.
<b>Less Than</b>	Determines if the field value is less than the value specified. This operation only works on numeric fields.
<b>Less Than Or Equal To</b>	Determines if the field value is less than or equal to the value specified. This operation only works on numeric fields.
<b>Longest</b>	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.
<b>Lowest</b>	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.
<b>Most Common</b>	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.
<b>Not Equal</b>	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><b>Note:</b> This option is not available if you select the operator Highest, Lowest, or Longest.</p> <p><b>Field</b> Choose this option if you want to compare another dataflow field's value to the field.</p> <p><b>String</b> 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 <b>Field type</b> field, select a dataflow field. If you selected String in the <b>Value type</b> field, type the value you want to use in the comparison.</p> <p><b>Note:</b> This option is not available if you select the operator Highest, Lowest, or Longest.</p>

4. Click **OK**.
5. 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.
6. Click the **Actions** node in the tree.
7. Click **Add Action**.
8. 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><b>Field</b> Choose this option if you want to copy a value from a field to the best of breed record.</p> <p><b>String</b> 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 <b>Source data</b> field.
Accumulate source data	<p>If the data in the <b>Source data</b> 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 8: Best of Breed Output**

Field Name	Format	Description / Valid Values
CollectionRecordType	String	Identifies the template and Best of Breed records in a collection of duplicate records. The possible values are: <b>Primary</b> The record is the selected template record in a collection. <b>Secondary</b> The record is not the selected template record in a collection. <b>BestOfBreed</b> The record is the newly created best of breed record in the collection.

## 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 9: Candidate Finder Database Options**

Option Name	Description / Valid Values
Finder type	Select Database.



Option Name	Description / Valid Values
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.  <b>Note:</b> 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 <a href="#">Defining the SQL Query</a> on page 177
Field Map tab	Choose field mapping settings as described in <a href="#">Mapping Database Columns to Stage Fields</a> on page 178.
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 10: Candidate Finder Options**

Option Name	Description / Valid Values
Finder type	Select Search Index.
Name	Select the appropriate index that was created using the <a href="#">Write to Search Index</a> on page 220 stage under the Advanced Matching deployed stages in Enterprise Designer.

Option Name	Description / Valid Values
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><b>Note:</b> If the maximum results is arbitrarily large, process those in batches, using the <b>Fetch Batch Size</b> field.</p>
Fetch Batch Size	<p>If the <b>Maximum results</b> 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 <b>Fetch Batch Size</b> is a value lesser than <b>Maximum results</b> and if the <b>Fetch Batch Size</b> is greater than <b>Maximum results</b>, the records are processed in a single batch.</p> <p><b>Note:</b> This field is applicable only to cluster supported search engine and not to the legacy search engine.</p>
Return match count	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.
Index search type	Determines the type of index search you want to conduct. Select <b>Simple search</b> .
Index Fields	Select the index field(s) you want to use for comparison in the simple search.
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
-------------	----------------------------

Specify which analyzer to use to tokenize the input string. One of these:

- 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 “Pitney Bowes Software” would be returned as three tokens: “Pitney”, “Bowes”, and “Software”.
- 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.
- StopWord—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. For example, the string “Pitney Bowes Software” would be returned as just one token “Pitney Bowes Software”.
- 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.

Option Name	Description / Valid Values
	<ul style="list-style-type: none"> <li>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.</li> <li>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.</li> </ul>
Output Fields tab	<p>Check the <b>Include</b> box to select which stored fields should be included in the output.</p> <p><b>Note:</b> 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 11: Candidate Finder Options**

Option Name	Description / Valid Values
Finder type	Select Search Index.
Name	Select the appropriate index that was created using the <a href="#">Write to Search Index</a> on page 220 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><b>Note:</b> If the maximum results is arbitrarily large, process those in batches, using the <b>Fetch Batch Size</b> field.</p>

Option Name	Description / Valid Values
Fetch Batch Size	<p>If the <b>Maximum results</b> 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 <b>Fetch Batch Size</b> is a value lesser than <b>Maximum results</b> and if the <b>Fetch Batch Size</b> is greater than <b>Maximum results</b>, the records are processed in a single batch.</p> <p><b>Note:</b> 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 <b>Sort</b> check-box, the desired index field from the <b>Sort by</b> drop-down list, and select <b>Ascending</b> or <b>Descending</b> from the <b>Order by</b> drop-down list.</p> <p><b>Note:</b> You can perform sorting only on <b>String Field Type</b> and <b>Keyword Analyzer</b>.</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 <b>Index Field</b>.</p>
Index search type	<p>Determines the type of index search you want to conduct. Select <b>Advanced search</b>.</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 "Pitney" and "Pitney Bowes" would be contained in a search index field of "Pitney Bowes Software Inc."</p>
Contains All	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.
Contains Any	Determines whether any of the alphanumeric words from the input field is contained in the search index field.

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 <b>Maximum edits</b> parameter to set a limit on the number of edits allowed to be considered a successful match:</p> <ul style="list-style-type: none"> <li>• 0—Allows for no deletions, insertions, or substitutions. The input field data and the search index field data must be identical.</li> <li>• 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".</li> <li>• 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".</li> </ul> <p>The Fuzzy search type is used for single-word searches only. Click <b>Ignore extra words</b> 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 "Pitney" and the input field says "Pitney Bowes", they would not be considered a match because of "Bowes". However, if you check this box, "Bowes" would be ignored and with "Pitney" 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 <b>Ignore extra words</b> 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 <b>Pattern string</b> 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 <b>Ignore extra words</b> 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"> <li>Define the input <b>First input field</b> and <b>Second input field</b> you want to search for in the index.</li> <li>Use the <b>Distance</b> 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.</li> </ul> <p>For example, you could successfully use this search type to look for First field "Spectrum" and Second field "Pitney" within ten words of each other in a search index field containing the sentence "Spectrum Technology Platform is a product of Pitney Bowes Software Inc."</p> <p>The Proximity search type is used for single-word searches only. Click <b>Ignore extra words</b> 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"> <li>Use the <b>Lower bound field</b> parameter to select the field to be used as the starting term.</li> <li>Use the <b>Upper bound field</b> parameter to select the field to be used as the ending term.</li> </ul> <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 <b>Ignore extra words</b> 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 <b>Wildcard characters</b>.</p> <p>Select the <b>Position</b> 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 <b>Ignore extra words</b> 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 <b>Index field name</b> and enter "5" here.</p> <p><b>Note:</b> 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><b>Note:</b> By default the query ignores the blank fields.</p>
Output Fields tab	<p>Check the <b>Include</b> box to select which stored fields should be included in the output.</p> <p><b>Note:</b> 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 `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.

## Output

**Table 12: 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:  <b>Y</b> - The record is suspect and has candidates <b>N</b> - The record is suspect and doesn't have candidates <b>D</b> - 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:  <b>Suspect</b> A suspect record is used as input to a query. <b>Candidate</b> 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 <b>Group by</b> 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="553 373 1429 520">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 <b>Override sort performance options</b> box then specify the values you want in these fields:</p> <p data-bbox="553 541 1429 926"><b>In memory record limit</b> 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="818 810 1429 926"><b>Note:</b> Be careful in environments where there are jobs running concurrently because increasing the <b>In memory record limit</b> setting increases the likelihood of running out of memory.</p> <p data-bbox="553 961 1429 1213"><b>Maximum number of temporary files</b> 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 data-bbox="732 1234 1414 1360" style="background-color: #f0f0f0; padding: 10px; margin: 10px 0;"> <math display="block">\frac{(NumberOfRecords \times 2)}{InMemoryRecordLimit} = NumberOfTempFiles</math> </div> <p data-bbox="732 1377 1429 1440">Note that the maximum number of temporary files cannot be more than 1,000.</p> <p data-bbox="553 1457 1429 1520"><b>Enable compression</b> Specifies that temporary files are compressed when they are written to disk.</p> <p data-bbox="639 1541 1429 1625"><b>Note:</b> The optimal sort performance settings depends on your server's hardware configuration. Nevertheless, the following equation generally produces good sort performance:</p> <div data-bbox="639 1646 1414 1745" style="background-color: #f0f0f0; padding: 10px; margin: 10px 0;"> <math display="block">(InMemoryRecordLimit \times MaxNumberOfTempFiles \div 2) \geq TotalNumberOfRecords</math> </div>

## 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: <ul style="list-style-type: none"> <li><b>Non-Numeric</b> Choose this option if the field contains non-numeric data (for example, string data).</li> <li><b>Numeric</b> Choose this option if the field contains numeric data (for example, double, float, and so on).</li> </ul>



Option	Description
Operator	Specifies the type of comparison you want to use to evaluate the field. One of the following:
<b>Contains</b>	Determines if the field contains the value specified. For example, "sailboat" contains the value "boat".
<b>Equal</b>	Determines if the field contains the exact value specified.
<b>Greater Than</b>	Determines if the field value is greater than the value specified. This operation only works on numeric fields.
<b>Greater Than Or Equal To</b>	Determines if the field value is greater than or equal to the value specified. This operation only works on numeric fields.
<b>Highest</b>	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.
<b>Is Empty</b>	Determines if the field contains no value.
<b>Is Not Empty</b>	Determines if the field contains any value.
<b>Less Than</b>	Determines if the field value is less than the value specified. This operation only works on numeric fields.
<b>Less Than Or Equal To</b>	Determines if the field value is less than or equal to the value specified. This operation only works on numeric fields.
<b>Longest</b>	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.
<b>Lowest</b>	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.
<b>Most Common</b>	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.
<b>Not Equal</b>	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><b>Note:</b> This option is not available if you select the operator Highest, Lowest, or Longest.</p> <p><b>Field</b> Choose this option if you want to compare another dataflow field's value to the field.</p> <p><b>String</b> 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 <b>Field type</b> field, select a dataflow field. If you selected String in the <b>Value type</b> field, type the value you want to use in the comparison.</p> <p><b>Note:</b> 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><b>Field</b> Choose this option if you want to copy a value from a field to the other records in the group.</p> <p><b>String</b> 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><b>Note:</b> 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>
Destination	<p>Specifies the field in the other records to which you want to copy the data specified in the <b>Source data</b> 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.</p>

**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 <b>Group by</b> 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="553 373 1429 520">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 <b>Override sort performance options</b> box then specify the values you want in these fields:</p> <p data-bbox="553 541 1429 926"> <b>In memory record limit</b> 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 data-bbox="816 808 1429 926"><b>Note:</b> Be careful in environments where there are jobs running concurrently because increasing the <b>In memory record limit</b> setting increases the likelihood of running out of memory.</p> </p> <p data-bbox="553 961 1429 1220"> <b>Maximum number of temporary files</b> 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:           <div data-bbox="732 1234 1414 1360" style="background-color: #f0f0f0; padding: 10px; margin: 10px 0;"> <math display="block">\frac{(NumberOfRecords \times 2)}{InMemoryRecordLimit} = NumberOfTempFiles</math> </div> <p data-bbox="732 1375 1429 1438">Note that the maximum number of temporary files cannot be more than 1,000.</p> </p> <p data-bbox="553 1459 1429 1522"> <b>Enable compression</b> Specifies that temporary files are compressed when they are written to disk.           <p data-bbox="638 1543 1429 1627"><b>Note:</b> The optimal sort performance settings depends on your server's hardware configuration. Nevertheless, the following equation generally produces good sort performance:</p> <div data-bbox="638 1648 1414 1745" style="background-color: #f0f0f0; padding: 10px; margin: 10px 0;"> <math display="block">(InMemoryRecordLimit \times MaxNumberOfTempFiles \div 2) \geq TotalNumberOfRecords</math> </div> </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><b>Note:</b> 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><b>Note:</b> 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.

Option	Description
Field Type	Specifies the type of data in the field. One of the following:
	<b>Non-Numeric</b> Choose this option if the field contains non-numeric data (for example, string data).
	<b>Numeric</b> 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:
<b>Contains</b>	Determines if the field contains the value specified. For example, "sailboat" contains the value "boat".
<b>Equal</b>	Determines if the field contains the exact value specified.
<b>Greater Than</b>	Determines if the field value is greater than the value specified. This operation only works on numeric fields.
<b>Greater Than Or Equal To</b>	Determines if the field value is greater than or equal to the value specified. This operation only works on numeric fields.
<b>Highest</b>	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.
<b>Is Empty</b>	Determines if the field contains no value.
<b>Is Not Empty</b>	Determines if the field contains any value.
<b>Less Than</b>	Determines if the field value is less than the value specified. This operation only works on numeric fields.
<b>Less Than Or Equal To</b>	Determines if the field value is less than or equal to the value specified. This operation only works on numeric fields.
<b>Longest</b>	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.
<b>Lowest</b>	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.
<b>Most Common</b>	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.
<b>Not Equal</b>	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><b>Note:</b> This option is not available if you select the operator Highest, Lowest, or Longest.</p> <p><b>Field</b> Choose this option if you want to compare another dataflow field's value to the field.</p> <p><b>String</b> 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 <b>Field type</b> field, select a dataflow field. If you selected String in the <b>Value type</b> field, type the value you want to use in the comparison.</p> <p><b>Note:</b> 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:

$$\text{NumberOfRecords} \times 2 \div \text{InMemoryRecordLimit} = \text{NumberOfTempFiles}$$

Note that 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. Nevertheless, the following equation generally produces good sort performance:

$$(\text{InMemoryRecordLimit} \times \text{MaxNumberOfTempFiles} \div 2) \geq \text{TotalNumberOfRecords}$$

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 210 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
<b>Compare suspect to all candidates</b>	<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 <b>Return Unique Candidates</b> box to return records within a match group from the candidate port that have been identified as unique records.</p>
<b>Stop comparing suspect against candidates after finding <i>n</i> duplicates</b>	<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 112.
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
<b>Collection Number</b>	<b>Record Type</b>
1	Unique
2	Unique
3	Unique
4	Duplicate/Suspect
4	Duplicate/Suspect
<b>Collection Number</b>	<b>Record Type</b>
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. If you are creating a new custom matching rule, see [Building a Match Rule](#) on page 68 for more information.
11. Click **Evaluate** to evaluate how a suspect record scored against candidate records. For more information, see [Interflow Match](#) on page 200.

## Output

**Table 13: 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: <ul style="list-style-type: none"> <li><b>Y</b>      The record is a suspect record and has duplicates.</li> <li><b>N</b>      The record is a suspect record and has no duplicates.</li> <li><b>D</b>      The record is a candidate record and is a duplicate of the suspect record.</li> <li><b>U</b>      The record is a candidate record but is not a duplicate of the suspect record.</li> </ul>
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: <ul style="list-style-type: none"> <li><b>suspect</b>      The original input record that was flagged as possibly having duplicate records.</li> <li><b>duplicate</b>      A record that is a duplicate of the input record.</li> <li><b>unique</b>      A record that has no duplicates.</li> </ul>
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.

**Note:** The Validate Address and Advanced Matching Module 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 Module 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:

$$\frac{(\text{NumberOfRecords} \times 2)}{\text{InMemoryRecordLimit}} = \text{NumberOfTempFiles}$$

Note that 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. Nevertheless, the following equation generally produces good sort performance:

$$\frac{(\text{InMemoryRecordLimit} \times \text{MaxNumberOfTempFiles} \div 2)}{\text{TotalNumberOfRecords}} \geq$$

- 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 210 for more information.
- 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.
- Click **Sliding Window** to enable this matching method. For more information about Sliding Window, see [Sliding Window Matching Method](#) on page 209
- Click **Generate Data for Analysis** to generate match results. For more information, see [Analyzing Match Results](#) on page 112.
- 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	<u>Record Type</u>

Option	Description
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. For information about modifying the other options, see [Building a Match Rule](#) on page 68.
11. Click **Evaluate** to evaluate how a suspect record scored against candidate records. For more information, see [Interflow Match](#) on page 200.

### 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 14: 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> <ul style="list-style-type: none"> <li><b>suspect</b>      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.</li> <li><b>duplicate</b>      A record that is a duplicate of the suspect record.</li> <li><b>unique</b>      A record that has no duplicates.</li> </ul>
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.

**Note:** The Validate Address and Advanced Matching Module 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 Module 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
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 15: Match Key Generator Options**

Option Name	Description / Valid Values
-------------	----------------------------

---

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

Option Name	Description / Valid Values
	Specifies the algorithm to use to generate the match key. One of the following:
<b>Consonant</b>	Returns specified fields with consonants removed.
<b>Double Metaphone</b>	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.
<b>Koeln</b>	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.
<b>MD5</b>	A message digest algorithm that produces a 128-bit hash value. This algorithm is commonly used to check data integrity.
<b>Metaphone</b>	Returns a Metaphone coded key of selected fields. Metaphone is an algorithm for coding words using their English pronunciation.
<b>Metaphone (Spanish)</b>	Returns a Metaphone coded key of selected fields for the Spanish language. This metaphone algorithm codes words using their Spanish pronunciation.
<b>Metaphone 3</b>	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.
<b>Nysiis</b>	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.
<b>Phonix</b>	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

Option Name	Description / Valid Values
	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.
	<b>Soundex</b> Returns a Soundex code of selected fields. Soundex produces a fixed-length code based on the English pronunciation of a word.
	<b>Substring</b> Returns a specified portion of the selected field.
Field name	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.
Start position	Specifies the starting position within the specified field. Not all algorithms allow you to specify a start position.
Length	Specifies the length of characters to include from the starting position. Not all algorithms allow you to specify a length.
Remove noise characters	Removes all non-numeric and non-alpha characters such as hyphens, white space, and other special characters from an input field.
Sort input	Sorts all characters in an input field or all terms in an input field in alphabetical order. <b>Characters</b> Sorts the characters values from an input field prior to creating a unique ID. <b>Terms</b> Sorts each term value from an input field prior to creating a unique ID.

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 16: 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 112.
4. For information about modifying the other options, see [Building a Match Rule](#) on page 68.
5. Click **Evaluate** to evaluate how a suspect record scored against candidate records. For more information, see [Interflow Match](#) on page 200.

## Output

**Table 17: 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><b>Y</b>      The record is a suspect record and has duplicates.</p> <p><b>N</b>      The record is a suspect record and has no duplicates.</p> <p><b>D</b>      The record is a candidate record and is a duplicate of the suspect record.</p> <p><b>U</b>      The record is a candidate record but is not a duplicate of the suspect record.</p>
MatchRecordType	<p>Identifies the type of match record in a collection. The possible values are:</p> <p><b>Suspect</b>      The original input record that was flagged as possibly having duplicate records.</p> <p><b>Duplicate</b>      A record that is a duplicate of the input record.</p> <p><b>Unique</b>      A record that has no duplicates.</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>

Field Name	Description / Valid Values
MatchInfo.<RuleName>.IsMatch	<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 <b>Return detailed match information</b> option in the <b>Transactional Match Options</b> screen.</p> <p><b>Note:</b> The value in this field is derived from the match state of the child nodes based on the defined parent rule configuration.</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. This information is displayed if you select the <b>Return detailed match information</b> option in the <b>Transactional Match Options</b> screen.</p> <p><b>Note:</b> The score in this field is derived from the scores of the child nodes on the basis of the defined parent rule configuration.</p>
MatchInfo.<RuleName>.<RuleNodeName>.IsMatch	<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 <b>Return detailed match information</b> option in the <b>Transactional Match Options</b> screen.</p>
MatchInfo.<RuleName>.<RuleNodeName>.Score	<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 <b>Return detailed match information</b> option in the <b>Transactional Match Options</b> 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 Module 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 Module 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 <b>Key field</b>	All the records are stored.  <b>Note:</b> 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 <b>overwritten</b> .
Update or Append	Duplicate records with same <b>Key field</b>	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 “Pitney Bowes Software” would be returned as three tokens: “pitney”, “bowes”, and “software”. For a comparison of Standard and Keyword analyzers, see [Standard and Keyword Analyzer](#) on page 224.
  - 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 “Pitney Bowes Software” would be returned as just one token “Pitney Bowes Software”. For a comparison of Standard and Keyword analyzers, see [Standard and Keyword Analyzer](#) on page 224.
  - 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 index(es) 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		
P O	1	Y				Suspect		
P O	1	D	P O	4	Test4	Candidate		
P O	1	D	P	1	Test 1	Candidate		

Search result with Keyword analyzer:

Point		Point 2						
SearchTerm	CandidateGroup	HasDuplicates	Index	IndexField	Term	TransactionRecordType		
P O	1	Y				Suspect		
P O	1	D	P O	4	Test4	Candidate		



# Business Steward Module

## Business Steward Module

The Business Steward Module is a set of features that allow you to identify and resolve exception records. Exception records are records that Spectrum™ Technology Platform could not confidently process and that require manual review by a data steward. Some examples of exceptions are:

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

The Business Steward Module provides a browser-based tool for manually reviewing exception records. Once exception records are manually corrected and approved, they can be reincorporated into your Spectrum™ Technology Platform data quality process.

Many of the actions that take place within the Business Steward Module are reflected in the audit log, a Management Console tool that records user activity. The following actions are included in the log:

- Adding exceptions in the Write Exceptions stage
- Deleting exceptions in the Read Exceptions stage and the Business Steward Portal Manage Exceptions page
- Assigning exceptions in the Business Steward Portal Exception Manager
- Retrieving exceptions in the Business Steward Portal Exception Editor
- Updating exceptions in the Business Steward Portal Exception Editor
- Revalidating exceptions in the Business Steward Portal Exception Editor

Read more about the audit log in the Spectrum™ Technology Platform Administration Guide for the WebUI.

### Components

The Business Steward Module consists of:

- **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.
- **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.

- **Exception Dashboard and Editor**—A browser-based tool that displays a dashboard of summary statistics and charts to help you understand the kinds of exceptions that occur in your data. It is also an editor that allows you to modify exception records and approve them for reprocessing.
- **Manage Exceptions**—A browser-based tool that enables you to review and manage exception record activity for all users.
- **Data Quality Performance**—A browser-based tool that provides information on trends within your exception records and enables you to identify key performance indicators and send notifications when certain conditions have been met.
- **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 Business Steward Portal).

In addition to setting conditions that determine if records require manual review, you can also configure Exception Monitor to send a notification to one or more email addresses when those conditions have been met a certain number of times.

For more information on exception processing, see [Business Steward Module](#) on page 225.

### 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.

## Options

### Conditions Tab

**Table 18: Exception Monitor Options**

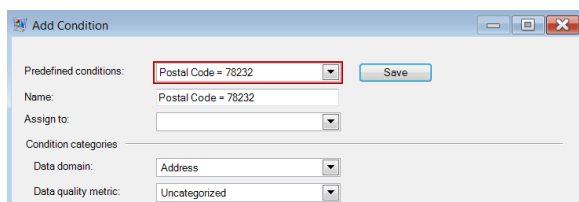
Option Name	Description
Stop evaluating when a condition is met	Specifies whether to continue evaluating a record against the remaining conditions once a condition is met. Enabling this option may improve performance because it potentially reduces the number of evaluations that the system has to perform. However, if not all conditions are evaluated you will lose some degree of completeness in the exception reports shown in the Business Steward 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.

### Adding or Modifying Conditions and Expressions

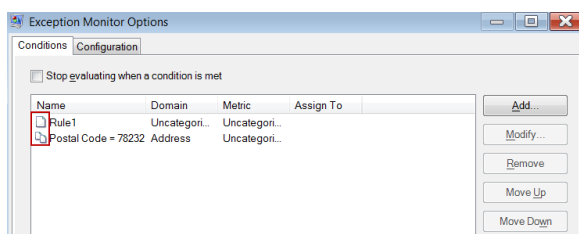
A condition defines the criteria used to determine if a record is an "exception" and needs to be routed for manual review. Typically this means that you want to define conditions that can consistently identify records that either failed automated processing earlier in the dataflow or that have a low degree of confidence and therefore should be reviewed manually.

The Exception Monitor stage enables you to create predefined conditions and custom conditions using the Add Condition dialog box. Predefined conditions are available to all dataflows, while custom conditions are available only to the dataflows for which they were created. The configuration process is almost identical for both types; however, to create a predefined condition you must save the condition by completing the fields and clicking **Save**, shown in the red box below.

After you have saved a custom condition, the **Predefined conditions** field changes to show the name of the condition rather than "<custom condition>".



After you have created predefined or custom conditions, they will appear on the Conditions tab of the Exception Monitor Options dialog box. As shown in the following image, 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.



1. In the Conditions tab of the **Exception Monitor Options** window, click **Add** to create a new condition, or **Modify** to edit an existing condition. Complete these fields:
  - **Predefined Conditions**—Select a predefined condition or retain "<custom condition>" in the dropdown to create a new condition.
  - **Name**—A name for the condition. The name can be anything you like. Since the condition name is displayed in the Business Steward 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"). Say "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 name is 100% identical.
  - **Assign to**—Select a user to whom the exception records meeting this condition should be assigned. If you do not make a selection in this field, the exception records will automatically be assigned to the user who ran the job.
  - **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.

2. You must add at least one expression to the condition. An expression is a logical statement that checks the value of a field. To add an expression, click **Add**. To modify an existing expression, click **Modify**. Complete these fields:
  - **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 231.
  - If other expressions are already defined for this condition, you can select an operator in the **Logical operator** field. One of the following:
    - **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 to create an expression with expression builder the following fields 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.
3. Click **Add** to add the expression. Click **Close** when you are done adding expressions.
4. Use the **Move Up** and **Move Down** buttons to change the order in which expressions are evaluated.
5. Click the **Notification** tab if you want Exception Monitor to send a message to one or more email addresses when this condition is met a specific number of times. That email will include a link to the failed records in the Exception Editor of the Business Steward Portal, where you can manually enter the correct data. If you do not wish to set up notifications, skip ahead to step 11. To stop receiving notifications at a particular email address, remove that address from the list of recipients in the **Send notification to** line of the Notification tab on the Modify Condition dialog box.

**Note:** Notifications must be set up in the Management Console before you can successfully use a notification from within Exception Monitor. See the *Administration Guide* for information on configuring notifications.

6. Enter the email address(es) to which the notification should be sent. Separate multiple addresses with commas, spaces, or semicolons.
7. Designate the point at which you want a notification to be sent. You can have it sent upon the first occurrence of the condition, or you can have it sent when the condition has been met a specific number of times. The maximum value is 1,000,000 occurrences.
8. Check the **Send reminder after** box if you want reminder messages sent to the designated email address(es) after the initial email.

9. Enter the number of days after the initial email that you want the reminder email to be sent.
10. Click **Remind daily** if you want reminder messages sent every day following the first reminder email.
11. Enter the **Subject** for the email notification.
12. If you want to save this condition for reuse as a predefined condition, click **Save**. If you modify an existing condition and click Save, you will be asked if you want to overwrite the existing condition; note that if you overwrite a predefined condition, those changes will take effect for all dataflows that use the condition.
13. When finished working with expressions, click **OK**.
14. Add or modify additional conditions as needed.
15. 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 information about, see [Configuration Tab](#) on page 233.
16. When finished, click **OK**.

#### *Removing a Condition or Expression*

- To remove a condition, open Exception Monitor, select the condition you want to remove, then click **Remove**. Note that when you remove a condition, all expressions in the condition are removed.
- To remove an expression, open the condition that contains the expression, select the expression, then click **Remove**.

#### *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](http://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.



## Configuration Tab

**Table 19: Exception Monitor Options**

Option Name	Description
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 <b>Stop job after reaching exception limit</b> 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 Business Steward Portal, but does not create exceptions for those records.
Return all records in exception's group	<p>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.</p> <p>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.</p> <p><b>Note:</b></p> <p>If the input data does not contain a field called "CollectionNumber" this option will be disabled.</p>
Group by	<p>If you selected <b>Return all records in exception's group</b>, choose the field by which to group the records.</p> <p><b>Note:</b> The "CollectionNumber" input field will not appear in this list because it is not a valid selection for the Group by feature.</p>

Option Name	Description
Revalidation service	Select the service you want to run when you revalidate records from this dataflow.
Action after revalidation	Specifies whether you want to reprocess records or approve records that have been revalidated.
Match exception records using match field	<p>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.</p> <p>The input records will be evaluated against the condition(s) and then matched against the existing exception records in the repository. If an input record passes the condition(s) and matches an exception record, that exception record will be removed from the repository. If an input record does not pass the condition(s) 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.</p>
Optimized for single records or small batches	<p>This option is activated when you select <b>Match fields</b>. When this option is unchecked (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.</p> <p>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.</p>
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 <b>Match exception records using match fields</b> box.

## Output

Exception Monitor returns records in two ports. One port contains records that do not meet any of the conditions defined in the Exception Monitor stage. The other port, the exception port, contains all 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 does not add or modify fields within a record.

## 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 [Business Steward Module](#) on page 225.)

**Note:** Once a record is read into a dataflow by Read Exceptions, it is deleted from the repository.

### Input

Read Exceptions reads in data from an exception repository. It does not take input from another stage in a dataflow.

**Note:** Only records marked as "approved" in the Business Steward Portal are read into the dataflow.

### Options

The Read Exceptions stage has the following options.

#### *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 Business Steward Portal and meet the filter criteria.

### Sort Tab

Use the Sort tab to sort the 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 on.
- **Remove:** Removes a sort field.

### Runtime Tab

- **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 Business Steward 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.

### Output

The Read Exceptions stage returns records from the exception repository that have been approved and that match the selection criteria specified in the Read Exception options. In addition to the records' fields, Read Exceptions returns these fields which describe the last modifications made to the record in the Business Steward Portal.

**Table 20: 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 business steward made to the record.

Field Name	Description
Exception.LastModifiedBy	The last user to modify the record in the Business Steward Portal
Exception.LastModifiedMilliseconds	The time that the record was last modified in the Business Steward 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 Business Steward 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

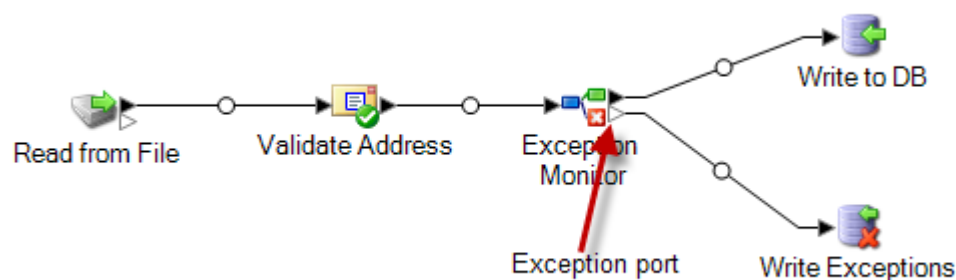
## Write Exceptions

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 Business Steward Portal.

**Note:** Exception records are not written to the Business Steward 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:



## 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 Business Steward 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 Business Steward Portal Exception 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 Business Steward Portal. Read more about best of breed records in the Business Steward Portal [here](#).

You may have input fields that you want in the repository but do not want to be viewable within the Business Steward 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 Business Steward Configuration tool or you can manually enter the name of the lookup. For more information on lookups, see [Lookups](#) on page 261.

**Note:** Lookups can be assigned only to fields whose type is string.

## Output

Write Exceptions does not return any output in the dataflow. It writes exception records to the exception repository.

# Business Steward Portal

## Business Steward Portal Introduction

### *What is the Business Steward Portal?*

The Business Steward Portal is a tool for reviewing, modifying, and approving records that failed automated processing or that were not processed with a sufficient level of confidence. Use the Business Steward Portal to manually enter correct or additional data in a record. 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 Portal to add information that was not in the original record.

The Business Steward Portal also provides summary charts that provide insight into the kinds of data that are triggering exception processing, including the data domain (name, addresses, spatial, and so on) as well as the data quality metric that the data is failing (completeness, accuracy, consistency, and so on).

In addition, the Business Steward Portal Manage Exception page enables you to review and manage exception record activity, including reassigning records from one user to another. Finally, the Business Steward Portal Data Quality page provides information regarding trends across dataflows and stages.

For more information on exception processing, see [Business Steward Module](#) on page 225.

### Accessing the Business Steward Portal

To open the Business Steward Portal, go to **Start > All Programs > Pitney Bowes > Spectrum Technology Platform > Server > Welcome Page** and select **Spectrum Data Quality**, then **Business Steward Portal**, and then click **Open the Business Steward Portal**.

Alternatively, you could follow these steps:

1. Open a web browser and go to `http://<servername>:8080/bsmportal`.

For example,

`http://myserver:8080/bsmportal`

Contact your Spectrum™ Technology Platform administrator if you do not know the server name and port.

2. Log in to the Spectrum™ Technology Platform. Contact your Spectrum™ Technology Platform administrator if you have trouble logging in.

### The Business Steward Portal Menu

The Business Steward Portal menu consists of four options and access to the help system, as shown below:

- **Dashboard**—View status of exceptions assigned to you. If you have view permissions, also view status of exceptions assigned to other users.
- **Editor**—Review, edit, and approve exception records for reprocessing.
- **Manage**—If you have view permissions, you can assign exceptions to yourself or others. If you have delete permissions, you can purge exception records from the repository.
- **Data Quality**—If you have view permissions, you can access this page to view statistical information and configure key performance indicators for exception records. If you do not have view permissions, you cannot access this page.
- **User Drop-Down**—Access the Business Steward Portal help system, all Spectrum Technology Platform documentation, and the Profile page, where you can set a language preference and specify a country to indicate culture preference, provide an email for notifications, and change your password.

### The Dashboard Page

The Exceptions Dashboard displays data that summarizes the status of exception records belonging to you and other users. (Note that you can only view others' data if you have modify permissions.) This includes the following:

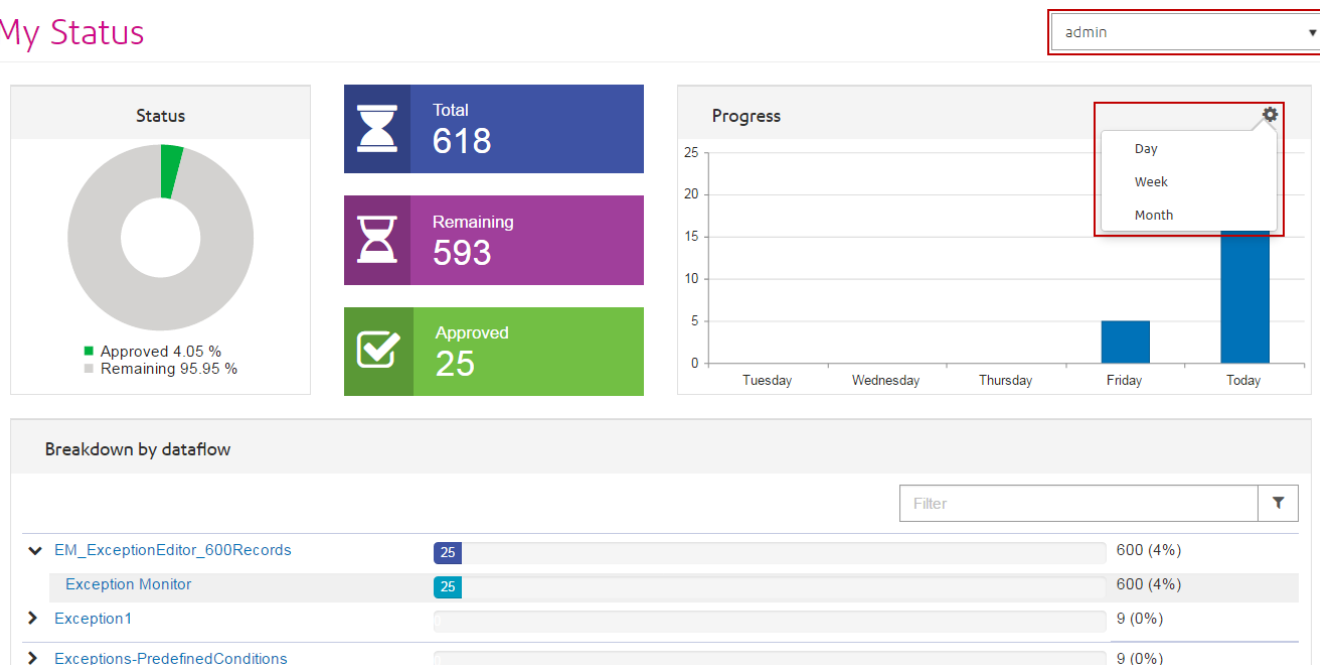
- The total number of exceptions
- The number and percentage of exceptions that have been approved versus those that remain unapproved
- The selected user's progress each day, week, or month
- Exception record approval progress by dataflow and stage(s) within the dataflow



**Note:** The progress charts for user and dataflow will not appear if you have turned off progress tracking on the Business Steward Settings page of Management Console.

- From the Dashboard tab, select the user whose exception activity you would like to view in the first drop-down box highlighted below. Note that only users who currently have exceptions assigned to them will appear in this list.

## My Status



- Click the **Settings** cog wheel (shown in the second drop-down box highlighted above) to select a duration of days, weeks, or months.
- Enter search criteria in the **Filter** field to narrow the list of dataflows or subset of dataflows. The search term must be included in the dataflow name.

## The Editor Page

The Exception Editor provides a means for you to manually review, modify, and approve exception records. The goal of a manual review is to determine which data is incorrect or missing and then update and approve it, particularly if Spectrum™ Technology Platform was unable to correct it as part of an automated dataflow process. You can then revalidate exception records for approval or reprocessing.

You can also use the Exception Editor to resolve duplicate exception records and use search tools to look up information that assists you in editing, approving, and rerunning records.

### Customizing Exception Editor Contents

There are multiple ways you can customize what is shown in the Exception Editor. You can use **Selection Options** to return records for a particular user, dataflow, job ID, and so on. You can use

the **Field Filter** tool to have records that meet certain criteria display, while records that don't meet the criteria are hidden.

If you use the **items per page** tool and then apply selection options or a field filter to narrow the list of exception records that are shown, don't forget that there may be multiple pages of results and not just what's shown on the initial screen. For example, let's say that you have set the items per page to 10 and then apply a field filter to return only the records that have a specific postal code. It may initially appear that only 10 records were returned, but there could be multiple pages of results since you set the limit of records shown at one time to 10.

### *Using Selection Options*

Select criteria in the drop-down fields of the selection options pane to narrow the records you see in the Exception grid.

<b>User</b>	Required. The ID of the user assigned to the dataflow. This information will be visible only if you have modify permissions.
<b>Dataflow name</b>	Required. The name of the dataflow that generated the exception records.
<b>Stage label</b>	Required. 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 exception record. The default label is "Exception Monitor".
<b>Job ID</b>	Required. A numeric identifier assigned to a job by the system. Each time a job runs it is assigned a new job ID.
<b>Status</b>	Optional. The approval status of the record.
<b>Date</b>	Optional. The date (and optionally time) that the dataflow ran. To enter time, type the time after the date.
<b>Data Domain</b>	Optional. 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.
<b>Quality Metrics</b>	Optional. 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.

### *Using the Field Filter*

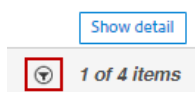
After you make selection options and the exception records are loaded, you can use field filtering to display only those records that you are interested in. By default, the Business Steward Portal only displays records from one Spectrum™ Technology Platform dataflow at a time. You can further filter the record list to show only the records that meet certain criteria within a particular field. Once a filter is created, it is automatically saved, and the next time you open the dataflow in the Exception Editor, the filter will be applied.

**Note:** You must apply selection options before using the field filter tool.

You can create filters for multiple fields, but you may create just one filter for each field. If a field already has a filter applied to it, the background of the arrow will be blue:

City	FirstName
KEENE	LAKSHMI
KEENE	LAKSHMI

Filters can only be created in the Tabular View. However, filters defined in the Tabular View are also reflected in the Form View. As with the Tabular View, an indicator is present near the bottom of the form to signify that a record has been filtered:



To filter the list of records:

1. Click the **Filter** button. You will see filter icons next to each column heading.
2. Click the filter icon for the field whose data you want to filter.
3. Select an operator that is appropriate for the field's data type, followed by a value.

**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.

<b>is greater than or equal to</b>	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.
<b>is less than</b>	Looks for records that have a numeric value that is less than the value you specify.
<b>is less than or equal to</b>	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.
<b>is after</b>	Looks for records that have a date or time value that is later than the value you specify.
<b>is after or equal to than</b>	Looks for records that have a date or time value that is equal to or later than the value you specify.
<b>is before</b>	Looks for records that have a date or time value that is earlier than the value you specify.
<b>is before or equal to than</b>	Looks for records that have a date or time value that is equal to or earlier than the value you specify.

4. Click the **Filter** button to apply the criteria. Only records whose data matches the criteria for that field will appear.
5. Click the filter icon again and click **Clear** to remove the filter. Or, click the **Filter** button to remove all filters; this action can be performed in either Tabular View or Form View.

### Viewing Records

You can view exception records in two different formats. The default view is the Tabular View, where you can load up to 100 exception records per page. You can scroll through the list and edit the records in any order. If you edit multiple records, you can save all the edits at one time; there is no need to save the edits for each individual record. If you use this view, you can specify how many records you want to see per page in the drop-down at the bottom of the screen.

The other method of viewing exception records is the Form View, where you view and edit one record at a time; you cannot edit multiple records at the same time in this view. Likewise, you must save the edits for each individual record; you cannot save the edits for multiple records at one time.

### Viewing Record Details

Regardless of which view you use, the Exception grid shows all the fields for a record as well as its approval 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 on the left end of a record. In the Form View, click **Show Detail**. These actions will open the **Detail** tab, which shows the following information.

<b>Job ID</b>	A numeric identifier assigned to a job by the system. Each time a job runs it is assigned a new job ID.
---------------	---

<b>Group By</b>	If the dataflow was configured to return all records in the exception 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.
<b>Exception Time</b>	The date and time when the Exception Monitor identified the record as an exception.
<b>Record Type</b>	The type of record you have selected. One of the following: <ul style="list-style-type: none"> <li>• <b>E</b>—Exception</li> <li>• <b>GE</b>—Group Exception</li> <li>• <b>N</b>—Best of Breed</li> </ul>
<b>Status</b>	The status of the record you have selected. One of the following: <ul style="list-style-type: none"> <li>• New</li> <li>• Resolved (possible only if the record has been edited at least once since its creation)</li> </ul>

Click the **Conditions** tab to view the following information:

<b>Condition</b>	The name of the condition that identified the record as an exception. Condition names are defined by the person who set up the dataflow.
<b>Domain</b>	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.
<b>Metric</b>	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.

Every time a record is changed, certain information is retained by the system, indicating who changed the record, when it was changed, the name of the user the record was assigned to, and any data that was entered in the comment field for that record. The record that appears in the Exception grid reflects the most recent changes and the most recent comment (if any); however, information on the **History** tab shows the following information for the entire life of the record, from when it was first added to the repository as an exception up to the point at which you are viewing the record:

<b>Version</b>	The revision number of the change.
<b>Last changed by</b>	The user who made the change.
<b>Assigned to</b>	The user to whom the exception record was assigned at the time of its revision.
<b>When</b>	The date and time that the change was saved.
<b>Comments</b>	The comments (if any) that were entered by the person who made the change.

### *Sorting Fields*

If you are using the Tabular View, you can sort the order of records shown based on a particular field by clicking anywhere in a column header. For instance, if you want to sort records in alphabetic

order by state, simply click the State column header. The first click will sort in ascending order, and a second click will sort in descending order; a third click will clear the sort order and return the records to their order prior to sorting.

### *Configuring Fields*

You can select which fields appear and change the order in which they appear by clicking the **Configure View** button (the cogwheel on the right side of the screen under the User Drop-Down) and making changes accordingly. These changes are saved on the server based on the user name and dataflow name; therefore, when you open the dataflow at a later time the configuration will still be applied. Similarly, changes you make here also affect what's shown when you edit exception records using the Form View. Use these features of Configure View to customize fields shown in the Exception Editor:

#### **Enabling Keyboard Navigation**

Check **Optimize tabular editor for keyboard navigation** to use the Tab key to navigate exception records in the Tabular view of the Exception Editor.

#### **Searching for Fields in Configure View**

Enter all or part of a field name in the Search box and the list of available fields will dynamically update. The search is case insensitive.

#### **Hiding Fields from View**

If you don't want to view every field from an exception record, click **Configure View** and deselect the fields you want to hide. The list shown will be in the same order as what you see in the Exceptions grid.

#### **Changing Field Order**

You can change the order in which fields are shown by dragging and dropping fields to put them in the desired order. However, you cannot rearrange fields when viewing search results. You must clear the search window and select "All" to resume the ability to reorder fields.


**Note:** If you are using the Tabular View, you can drag and drop column headings directly from the Exception Editor to change the order in which fields are shown; you do not need to do this from within Configure View.

#### **Configure the Form View Layout**

You can override the default layout of the fields in the Form View not only by selecting which columns should appear, but also by specifying their column width. The only rule in doing so is that the column widths for all visible fields must add up to a total of 12. For example, if you wanted four fields to be placed side-by-side on a single row, you could give each field a column width of "3." Likewise, if you wanted to include just two fields, you could enter "6" as the column width for each field. Alternatively, if you wanted to include five address fields on a row, you would need to specify a number of columns for each of those fields such that they add up to 12. The first image below depicts how you might configure five fields; the second image shows you an example of how that configuration renders.


**Note:** When you override the default layout for any field or combination of fields, those fields will no longer be responsive to the size of the window; however, fields with default values will still be responsive based on the size of browser window.

Configure view




Filter   All Visible Hidden

Field name	Form columns
<input checked="" type="checkbox"/> Approved	1
<input type="checkbox"/> Type	Default
<input checked="" type="checkbox"/> Address1	4
<input type="checkbox"/> FirmName	Default
<input checked="" type="checkbox"/> City	4
<input checked="" type="checkbox"/> StateProvince	1
<input checked="" type="checkbox"/> PostalCode	2

Optimize tabular editor for keyboard navigation

Reset to default 

OK Cancel

   Filter Table Form Settings

Approved

Address1	City	StateProvince	PostalCode
1594 Spring St.	Smyrna		30080

Show detail ▾

« < > » 1 of 11 items

### Editing Exception Records

The purpose of editing an exception record is to correct or augment and approve the record so that it can be processed successfully. Editing an exception 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 and approve them directly in the Exception Editor. You can edit one or more records at a time 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 exception 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 exception records with editable fields. Also, you cannot edit a field to contain a value that does not match the data type. For example, you cannot edit a field with a numeric data type to contain letters.

### *Tabular View*

To edit a field for a single record in the Tabular View:

1. Click the field you want to edit and change the field value accordingly. Read-only fields will be grayed out. Right-click the field to access cut, copy, and paste options. When you have edited a field, you will notice a green triangle appear in the upper-left corner of that field. This is a visual cue to remind you that the value of the field has been changed but is not yet saved.
2. Check the **Approved** box for the modified record. This will mark the record as ready to be processed by Spectrum™ Technology Platform.
3. If you need to undo a change you made, select the record you want to undo and click the **Undo changes** button.
4. Click the **Save** button when you are finished editing records.

To edit a field for multiple records in the Tabular View:

1. While pressing the Ctrl or Shift key, click the field you want to change for all the selected records. For instance, if you want to change all instances of "L.A." to say "Los Angeles," press the Ctrl key and click within the City field for one of the records. Then, while still pressing Ctrl or Shift, click within that same field for the other records you want to change.
2. Change the field values accordingly. You are able to edit these fields, but be aware that changes you make here will apply to all selected records, even though previously the values for those fields varied. Likewise, if you clear the data for a field when editing multiple records, it will be cleared for all selected records.
3. Check the **Approved** box for the modified records. This will mark the record as ready to be processed by Spectrum™ Technology Platform. Alternatively, you can click the **Approve all** button; this will result in every simple exception record showing in the Editor to be approved. (The Approval all feature does not apply to matching exception records.)
4. If you need to undo a change you made, select the records you want to undo and click the **Undo changes** button.
5. Click the **Save** button when you are finished editing records. The records' changes are saved to the exception repository and the view is refreshed. If you have not defined a revalidation workflow, the list of exceptions is reloaded. However, one or more of the edited records may not be in the refreshed list because they no longer match the search/filter criteria. If you *have* defined a revalidation workflow, the edited records may not be in the refreshed list if they are now valid and have been purged from the repository.
6. Use the navigation buttons at the bottom of the screen to go to previous or next page of exception records; you can also use these buttons to go directly to the first or last exception record.

### *Form View*

To edit records using the Form View, follow these steps:



1. Click **Form View**. The first record in the set will appear.
2. Click the field you want to edit, and change the field value accordingly. Read-only fields will be grayed out. Right-click the field to access cut, copy, and paste options. When you have edited a field, you will notice the outline of that field turn green. This is a visual cue to remind you that the value of the field has been changed but is not yet saved.
3. You can add comments about your changes in the **Comments** column. Comments are visible to other users and can be used to help keep track of the changes made to the record.
4. When you are confident that you have made the necessary changes to make the record valid, check the **Approved** box. This will mark the record as ready to be processed by Spectrum™ Technology Platform.
5. If you need to undo a change you made, click the **Undo changes** button.
6. Click **Save**. The record's changes are saved to the exception repository and the view is refreshed. 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.
7. Use the navigation buttons at the bottom of the screen to go to previous or next exception record; you can also use these buttons to go directly to the first or last exception record.

### *Resolving Duplicate Records*

Duplicate resolution exceptions occur when 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.

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 Business Steward 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. Select the record you want to work on then click **Resolve Duplicates**.

The **Duplicate Resolution** view shows duplicate records. The records are grouped into collections or candidate groups that contain these match record types:

<b>suspect</b>	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.
<b>duplicate</b>	A record that is a duplicate of the suspect record.
<b>unique</b>	A record that has no duplicates.

You can determine a record's type by looking at the MatchRecordType column.

2. If necessary, correct individual records as needed. For more information, see [Editing Exception Records](#) on page 247. Alternatively, you can drag and drop records across groups.
3. In the CollectionNumber or CandidateGroup field, enter the number of the group that you want to move the record into. 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 records, check the **Approved** box. This signals that the record is ready to be re-processed by Spectrum™ Technology Platform.
5. To save your changes, click the **Save** button.

#### *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 exception records.

1. Select the record you want to work on then click **Resolve Duplicates**.

The **Duplicate Resolution** view shows duplicate records. The records are grouped into collections or candidate groups that contain these match record types:

<b>suspect</b>	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.
<b>duplicate</b>	A record that is a duplicate of the suspect record.
<b>unique</b>	A record that has no duplicates.

You can determine a record's type by looking at the MatchRecordType column.

2. If necessary, correct individual records as needed. For more information, see [Editing Exception Records](#) on page 247.
3. Select a record that you want to put in the new collection then click **New Collection**. 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 records, check the **Approved** box. This signals that the record is ready to be re-processed by Spectrum™ Technology Platform.
6. To save your changes, click the **Save** button.

#### *Making a Record Unique*

To change a record from a duplicate to a unique:

1. In the MatchRecordType field, enter "Unique".
2. When you are done modifying records, check the **Approved** box. This signals that the record is ready to be re-processed by Spectrum™ Technology Platform.
3. To save your changes, click the **Save** button.

#### *Fields Automatically Adjusted During Duplicate Resolution*

When you modify records in the Business Steward Portal's duplicate resolution view, some fields are automatically adjusted to reflect the record's new disposition.

**Table 21: Records Processed by Interflow or Intraflow Match**

Action	Values Automatically Applied to Fields
Moving a record from one collection to another	<p><b>If you move a record into a collection of duplicates:</b></p> <ul style="list-style-type: none"> <li>• MatchRecordType: Duplicate</li> <li>• MatchScore: 100</li> <li>• HasDuplicates: D (This field is only present if the dataflow contained an Interflow Match stage.)</li> </ul> <p><b>If you move a duplicate record into the collection of unique records (collection 0):</b></p> <ul style="list-style-type: none"> <li>• MatchRecordType: Unique</li> <li>• MatchScore: No change</li> <li>• HasDuplicates: U (This field is only present if the dataflow contained an Interflow Match stage.)</li> </ul> <p><b>If you move a suspect record into the collection of unique records (collection 0):</b></p> <ul style="list-style-type: none"> <li>• MatchRecordType: Unique</li> <li>• MatchScore: 0</li> <li>• HasDuplicates: N (This field is only present if the dataflow contained an Interflow Match stage.)</li> </ul>
Creating a new collection	<ul style="list-style-type: none"> <li>• MatchRecordType: Suspect</li> <li>• MatchScore: No value</li> <li>• HasDuplicates: Y (This field is only present if the dataflow contained an Interflow Match stage.)</li> </ul> <p><b>Note:</b> 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 22: Records Processed by Transactional Match**

Action	Values Automatically Applied to Fields
Change MatchRecordType to Duplicate	<ul style="list-style-type: none"> <li>• HasDuplicates: D</li> <li>• MatchScore: 100</li> </ul>
Change MatchRecordType to Unique	<ul style="list-style-type: none"> <li>• HasDuplicates: U</li> <li>• MatchScore: unchanged</li> </ul>

Action	Values Automatically Applied to Fields
--------	--

Change HasDuplications to D	<ul style="list-style-type: none"> <li>MatchRecordType: Duplicate</li> <li>MatchScore: 100</li> </ul>
-----------------------------	---

Change HasDuplications to U	<ul style="list-style-type: none"> <li>MatchRecordType: Unique</li> <li>MatchScore: unchanged</li> </ul>
-----------------------------	--

Change HasDuplications to Y	<ul style="list-style-type: none"> <li>MatchRecordType: Suspect</li> <li>MatchScore: blank</li> </ul>
-----------------------------	---


Change HasDuplications to N	<ul style="list-style-type: none"> <li>MatchRecordType: Suspect</li> <li>MatchScore: blank</li> </ul>
-----------------------------	---

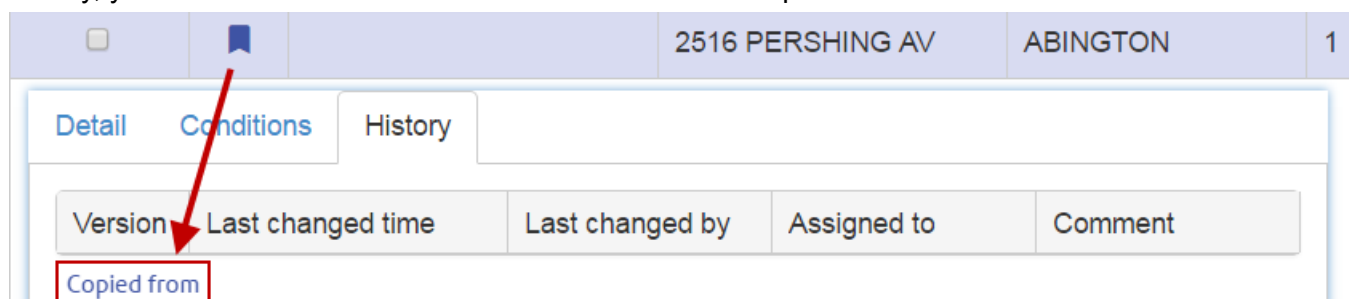
### Creating a Best of Breed Record

Follow these steps to create a best of breed record.

**Note:** 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 approve all the records within that group.

- To create a best of breed record within a collection, select a record in that 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 Exception 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 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. Note that 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.

See [Write Exceptions](#) options for more information on best of breed functionality.

### *Using Search Tools*

The Business Steward Portal Exception Editor provides search tools to assist you in looking up information that may help you edit and approve exception records and rerun them successfully. The tools include Bing Maps and the 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 Business Steward 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 Bing Maps*

Bing Maps displays the location of an address on a map and provides controls that allow you to zoom and pan the map. In addition, you can click the map to obtain addresses.

**Note:** The Bing Maps search tool is provided by Microsoft; you must be connected to the Internet to use this service.

1. 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 **Bing Map**.
4. Select the fields you want to use in your search from the **Field Name** column. For example, if you want to search for an address on a map, you might choose AddressLine1 and City. If you want to view just the city on a map, you might select just City and StateProvince.
5. Select "Road" for a traditional map view, "Bird's Eye" for an aerial view, and "Automatic" to let Bing Maps decide which is more appropriate.

6. Click **Go**. The results, including latitude and longitude, are displayed in the **Results** box and on the map in the form of a pin. Click the Rotate buttons to shift the perspective 90 degrees. Click the arrows on the compass to shift the focus incrementally in the selected direction. Use the Zoom In and Zoom Out buttons to focus more or less closely. You can also move the pin by dragging it to a new location; the map information will update dynamically.
7. To obtain the address of other buildings, click the map. Switching to Bird's Eye view may be helpful when finding buildings.

After completing an initial map search, you can click another exception record and the Go button, and the map will update accordingly.

#### *Using Spectrum Service Search Tools*

Pitney Bowes 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 Exception 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 BSM Portal Service drop-down.

1. Select the record that contains the data you want to look up.
2. Below the Exception Editor, click **Search Tools**.
3. In the **Service** field, select 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 your exception record doesn't 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 Business Steward Portal remembers the maps you create from service fields to exception fields as long as you are mapping exception records with the same field names, stage names, and dataflow names. For instance, if your exception 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. Repeat step 4 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 exception 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 exception 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 exception 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 enables a user with view and modify permissions to review and manage exception record activity for all assignees. It also provides the ability to reassign exception records from one user to another. If you have delete permissions, you can delete the entire group of exception records from the system based on dataflow name and job ID.

### Reviewing Exception Record Activity

The **Status** section of the Manage Exceptions page shows exception 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 of the screen. "Assignments" provides the number of exception 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 non-exception records that exist within groups that also contain exception records. Note that non-exception data will appear only if there are non-exception records within a group also containing exception records 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 each of 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 row. The list will dynamically auto-populate with dataflows or assignees that match the letters you type.

The **Progress** section of this page shows the cumulative progress for all users or all dataflows whose individual statuses are shown on the Dashboard. Select "Assignments" at the top of the screen to view user progress; select "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 Exception Records

The **Assignment** section of the Manage Exceptions page (under **Maintenance**) enables you to reassign exception records from one user to another.

1. Select a user whose exceptions you want to assign to another user in the **User** field.
2. To reassign all exception records belonging to a user, skip to Step 4. To reassign a portion of a user's exception records, complete one or more of these fields:
  - Dataflow name—The name of the dataflow producing the exception records.
  - Stage label—The name of the stage producing the exception records.
  - Job ID—The ID assigned to the job containing the exception records.
  - 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 exception records were created.
  - To date/time—The end date in a range of dates in which the exception records were created.
  - Approval status—Whether or not the exception records have been approved.
3. After making selections in the **User**, **Dataflow name**, and **Stage label** fields (at minimum), you can further refine the filter based on exception field values.
  - a) Click the add field filter icon.
  - b) In the **Field Name** column, select the field you want to filter on.
  - c) In the **Operator** column, select one of the following:
 

<b>is equal to</b>	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."
<b>is not equal to</b>	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."
<b>is greater than</b>	Looks for records that have a numeric value that is greater than the value you specify.
<b>is greater than or equal to</b>	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.
<b>is less than</b>	Looks for records that have a numeric value that is less than the value you specify.
<b>is less than or equal to</b>	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 **Field 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:

Field name	Operator	Value
State	Is equal to	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

- Click **Reassign**.
- 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.

6. Select another user in the **Reassign** dropdown.
7. Click **Confirm**.

### *Deleting Exception Records*

Occasionally you may want to delete exception records from the 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 of the Manage Page (under **Maintenance**) enables you to do this.

You must make selections from both the **Dataflow name** and **Job ID** fields before clicking **Remove**. However, you can select "All" from the Job ID field to remove exception records from every job run by the selected dataflow. Click **Remove data quality report data** to remove all performance data. If this option is not selected, the job's exception records will be removed from the repository but the performance data will still appear on the Data Quality page.

### **The Data Quality Page**

The Business Steward Portal Data Quality page provides information on trends within your exception records.

**Note:** This page will be deactivated if you turn off data quality reporting on the Business Steward 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 that 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.

## Business Steward Settings

### Introduction

#### *Business Steward Settings Introduction*

Business Steward Settings provides the following tools for users with write permissions:

- **Lookups**—Provides a method of correcting exception records using a specific set of data for those corrections.
- **Domains**—Enables you to specify the kind of data being evaluated.
- **Metrics**—Enables you to specify the way in which data is measured.
- **Notifications**—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**—Sets preferences for tracking pass/fail conditions and KPIs.
- **Search Tools Services**—Sets preferences for search tools in the Business Steward Portal Exception Editor.
- **Options**—Sets preferences for audit logs and progress tracking.

#### *Accessing Business Steward Settings*

To access Business Steward Settings:

1. In a web browser go 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** button.
4. Select **Business Steward Settings**.

## 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 Exception Editor. 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 Exception 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 in the Exception Editor of the Business Steward Portal 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 in the Exception Editor.

**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. Click the **Add lookup** button.

2. Enter a name for the new lookup in the text box.
3. To manually add a value/label pair:
  - Click the **Add lookup value** button.
  - Enter a value and/or a label for the lookup.

To use the **Add many** function:

- Click the **Add many** button to open the dialog box.
- Select the columns and separator accordingly. If you are pasting data from Microsoft Excel, use the Tab delimiter. If you define the wrong delimiter, the tool will import the entire line as the value or the label (whichever is designated as the first column).
- Type the values, separators, and/or labels for all the contents or paste the contents from another application.

After all values or value/label pairs have been added you can sort them in ascending or descending order on either the Value or the Label column. Note that once you have sorted the list, you cannot unsort it (except to sort it in reverse order or sort on the other column).

4. Repeat step 3 as necessary.
5. 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 the Business Steward Portal.

1. In the Exception Editor, 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 that 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".
4. Repeat step 3 for each problematic record.
5. Save the changed exceptions.

### **Modifying or Deleting Lookups**

1. Open the Lookups page.
2. Check the box next to the appropriate lookup.

3. Click the **Edit lookup** button, modify the lookup as necessary and click **Save**.

or

Click the **Delete lookup** 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 Business Steward Configuration 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 Business Steward Configuration 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. That email will include a link to the failed records in the Exception Editor of the Business Steward Portal, where you can manually enter the correct data. To stop receiving notifications at a particular email address, remove that address from the list of recipients in the Send notification to line of the Edit domain page.

**Note:** Notifications must be set up in the Management Console before you can successfully use a notification from within Business Steward Configuration. See the Administration Guide for information on configuring notifications.

1. In Business Steward Configuration, open the **Domains** page or the **Metrics** page.
2. Select the domain or metric that you want to add a notification for and click the **Edit item** button.
3. Select user name(s) from the drop-down list (as configured in Management Console) or enter new email addresses that the notifications should be sent to.
4. Select the number of exception records that should trigger a notification.
5. Enter the text that should be sent as the subject of the notification.
6. Enter the message that should appear in the body of the notification. You can use variables in the message to relay important information about the exception(s), including the following:
  - `#{jobID}`—The ID number of the job that produced the exception record(s).
  - `#{jobName}`—The name of the job that produced the exception record(s).



- `${userName}`—The name of the user whose job that produced the exception record(s).
  - `${stageLabel}`—The name of the dataflow stage that produced the exception record(s).
  - `${link}`—A link to the Editor page in the Business Steward Portal, showing records for a particular dataflow.
7. Check the **Send reminder** box if you want to send a reminder notification, and select the number of days that should pass before the reminder is sent.
  8. Select user name(s) from the drop-down list (as configured in Management Console) or enter an email address in the **Send reminder to** field to send notifications to email addresses other than the one(s) specified in step 3.
  9. Enter the text that should be sent as the subject of the reminder notification.
  10. Enter the message that should appear in the body of the reminder notification. The reminder uses an additional variable:
    - `${Count}`—The number of exceptions for the specified dataflow or stage that have yet to be resolved.
  11. Check the **Remind daily** box if you want a reminder notification to be sent every day until the exceptions are resolved.

## Data Quality Reporting

1. Click **Data quality reporting** to track pass or fail conditions in the Exception Monitor stage. If you turn off this option, the Business Steward Portal Data Quality page 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 Business Steward Module 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 Business Steward Module 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 the Business Steward Module 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

1. Select which search tool services you want to be available in the Business Steward Portal Exception 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

1. Click **Audit exception events** to have the Business Steward Module 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 the Business Steward Portal. If you turn this option off, the progress charts on the Business Steward Portal Dashboard will not appear.

# Data Normalization Module

## Data Normalization Module

The Data Normalization Module examines terms in a record and determines if the term is in the preferred form.

### Components

The Data Normalization Module consists of:

- **Advanced Transformer**—This stage scans and splits strings of data into multiple fields, placing the extracted and non extracted data into an existing field or a new 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 includes changing full words to abbreviations, changing abbreviations to full words, changing nick names to full names or misspellings to corrected spellings.
- **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.

Option	Description
Extract using	<p>Select <b>Table Data</b> or <b>Regular Expressions</b>.</p> <p>Select <b>Table Data</b> if you want to scan and split using the XML tables located at &lt;Drive&gt;:\Program Files\Fitney Boves\Spectrum\server\modules\advancedtransformer\data. See Table Data Options below for more information about each option.</p> <p>Select <b>Regular Expressions</b> if you want to scan and split using regular expressions. Regular expressions provide many additional options for splitting data. You can use the pre-packaged 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 an "123 Main St." would go in another. See Regular Expression options below for more information about each option.</p>
<b>Table Data Options</b>	
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>
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"> <li>• Smith</li> <li>• ,</li> <li>• John</li> </ul> <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>

Option	Description						
Table	<p>Specifies the table that contains the terms on which to base the splitting of the field. For a list of tables, see <a href="#">Advanced Transformer Tables</a> on page 151. For information about creating or modifying tables, see <a href="#">Introduction to Lookup Tables</a> on page 151.</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 <b>City</b>.</p> <p>For multiple word searches, the search stops at the first occurrence of a match.</p> <p><b>Note:</b> Selecting this option may adversely affect performance.</p>						
Extract	<p>Specifies the type of extraction to perform. Select from one of these:</p> <table border="0"> <tr> <td data-bbox="552 976 828 1018"><b>Extract term</b></td> <td data-bbox="844 976 1429 1018">Extracts the term identified by the selected table.</td> </tr> <tr> <td data-bbox="552 1029 828 1092"><b>Extract N words to the right of the term</b></td> <td data-bbox="844 1029 1429 1155">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.</td> </tr> <tr> <td data-bbox="552 1165 828 1228"><b>Extract N words to the left of the term</b></td> <td data-bbox="844 1165 1429 1291">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.</td> </tr> </table> <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>	<b>Extract term</b>	Extracts the term identified by the selected table.	<b>Extract N words to the right of the term</b>	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.	<b>Extract N words to the left of the term</b>	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.
<b>Extract term</b>	Extracts the term identified by the selected table.						
<b>Extract N words to the right of the term</b>	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.						
<b>Extract N words to the left of the term</b>	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.						

---

## Regular Expressions Options

---

Option	Description
Regular Expressions	<p>Select a pre-packaged 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 <b>java.util.regex</b>, enabling the use of regular expressions. For more information, go to: <a href="http://java.sun.com/docs/books/tutorial/essential/regex/index.html">java.sun.com/docs/books/tutorial/essential/regex/index.html</a>.</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 <b>Populate Group</b> to extract any RegEx groups and place the complete expression, as well as any RegEx groups found, into the Groups list.
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:  <code>(1[012]{1,2}0?[1-9])[-./]([12][0-9] 3[01]{1,2}0?[1-9])[-./]([0-9]{4})</code>. 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 (i.e. "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 Name 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.

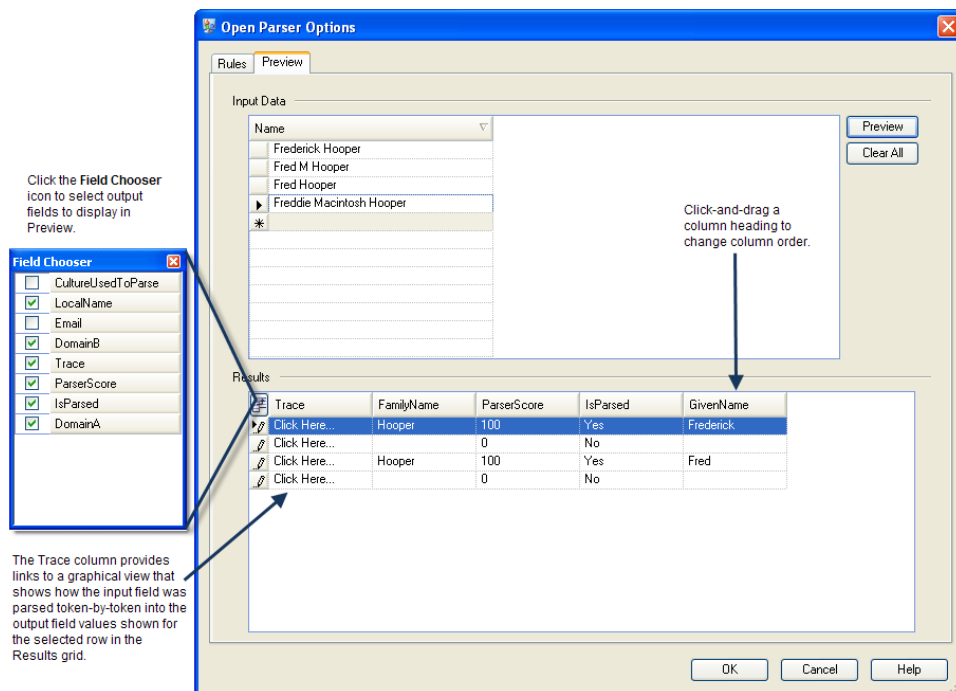
### Rules Tab

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 <a href="#">Defining a Culture-Specific Parsing Grammar</a> on page 10.</p> <p>If you choose this option you will also see these options:</p> <p><b>Domain</b> Specifies the parsing grammar to use.</p> <p><b>Cultures</b> Specifies the language or culture of the data you want to parse. Click the <b>Add</b> 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 <b>Move Up</b> and <b>Move Down</b> buttons. For more information about cultures, see <a href="#">Defining a Culture-Specific Parsing Grammar</a> on page 10.</p> <p><b>Return multiple parsed records</b> 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 input data's language or domain. 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><b>Note:</b> You can also define domain-independent grammar at runtime. Click <a href="#">here</a> for more information.</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 274. 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.

Field Name	Description / Valid Values
CultureCode	The culture codes contained in the input data. For a complete list of supported culture codes, see <a href="#">Assigning a Parsing Culture to a Record</a> on page 12.
CultureUsedtoParseSelect a match results in the Match Results List and then click <b>Remove</b> .	The culture code value used to parse each output record. This value is based on matches to a culture-specific parsing grammar.
IsParsed	Indicates whether an output record was parsed. The possible values are Yes or No.
ParserScoreSelect a match results in the Match Results List and then click <b>Remove</b> .	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 <a href="#">Scoring</a> .
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 <a href="#">Table Lookup Tables</a> on page 154.

## 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 data-bbox="535 367 1429 409">Specifies the type of action to take on the source field. One of the following:</p> <p data-bbox="535 420 1429 546"><b>Standardize</b> 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 data-bbox="535 556 1429 693"><b>Identify</b> Flags the record as containing a term that can be standardized, but performs no action on the data in the field. The output field <b>StandardizedTermIdentified</b> is added to the record with a value of Yes if the field can be standardized and No if it cannot.</p> <p data-bbox="535 703 1429 840"><b>Categorize</b> 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 data-bbox="535 913 1429 976">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 data-bbox="535 987 1429 1354"><b>Complete field</b> Treats the entire field as one term, resulting in the following:</p> <ul data-bbox="682 1039 1429 1354" style="list-style-type: none"> <li data-bbox="682 1039 1429 1165">• If you selected the action <b>Standardize</b>, 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".</li> <li data-bbox="682 1165 1429 1260">• If you selected the action <b>Identify</b>, Table Lookup treats the entire field as one string and flags the record if the string as a whole can be standardized.</li> <li data-bbox="682 1260 1429 1354">• If you selected the action <b>Categorize</b>, Table Lookup treats the entire field as one string and flags the record if the string as a whole can be categorized.</li> </ul> <p data-bbox="535 1365 1429 1873"><b>Individual terms within field</b> Treats each word in the field as its own term, resulting in the following:</p> <ul data-bbox="682 1417 1429 1873" style="list-style-type: none"> <li data-bbox="682 1417 1429 1543">• If you selected the action <b>Standardize</b>, 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."</li> <li data-bbox="682 1543 1429 1638">• If you selected the action <b>Identify</b>, Table Lookup parses the field and flags the record if any single term within the field can be standardized.</li> <li data-bbox="682 1638 1429 1873">• If you selected the action <b>Categorize</b>, Unlike <b>Standardize</b>, <b>Categorize</b> does not copy the source term if there isn't a table match. If none of the source terms match, <b>Categorize</b> uses the default value specified. Unlike <b>Standardize</b>, <b>Categorize</b> only returns that table value and nothing from <b>Source</b>. If none of the source terms match, <b>Categorize</b> uses the default value specified.</li> </ul>

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 <b>Destination</b> field as you did in the <b>Source</b> field. You can also create a new field by typing the name of the field you want to create.</p> <p>The <b>Destination</b> field is not available if you select the action <b>Identify</b>.</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 <a href="#">Table Lookup Tables</a> on page 154. For information about creating or modifying tables, see <a href="#">Introduction to Lookup Tables</a> on page 151.</p>				
Lookup multiple word terms	<p>Enables multiple word searches within a given string. For example:</p> <p>Input String: "Major General John Smith"            Business Rule: Identify "Major General" in a string based on a table that contains the entry            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 <b>On</b> is set to <b>Complete field</b>.</p> <p><b>Note:</b> 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><b>Source's value</b></td> <td>Put the value from the source field into the destination field.</td> </tr> <tr> <td><b>Other</b></td> <td>Put a specific value into the destination field.</td> </tr> </table>	<b>Source's value</b>	Put the value from the source field into the destination field.	<b>Other</b>	Put a specific value into the destination field.
<b>Source's value</b>	Put the value from the source field into the destination field.				
<b>Other</b>	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 <b>Complete field</b> or <b>Individual terms in field</b> options.
	<b>Yes</b> The record contains a term that can be standardized.
	<b>No</b> 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	
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](http://www.eki.ee/wgrs).

<b>Arabic</b>	The script used by several Asian and African languages, including Arabic, Persian, and Urdu.
<b>Cyrillic</b>	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.
<b>Devanagari</b>	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.
<b>Greek</b>	The script used by the Greek language. This script belongs to the Hellenic branch of the Indo-European language family.
<b>Gujarati</b>	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.
<b>Gurmukhi</b>	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.
<b>Hangul</b>	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 <a href="#">The National Institute of the Korean Language</a> .
<b>Han</b>	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.
<b>Traditional/Simplified Chinese</b>	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
<b>Kannada</b>	The script used by several South Indian languages, such as Konkani. This script is a descendent of Brahmi script of ancient India.
<b>Katakana and Hiragana</b>	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.

<b>Half width/Full width</b>	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.
<b>Latin</b>	The script used by most languages of Europe, such as English. It was originally used by the ancient Romans to write the Latin language.
<b>Malayalam</b>	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.
<b>Oriya</b>	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.
<b>Tamil</b>	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.
<b>Telugu</b>	The script used by several languages of South India. This script is a descendent of Brahmi script of ancient India.
<b>Thai</b>	The script used by Thai language. This script is influenced by the Brahmi script of ancient India and the Khmer alphabets.

Transliterator is part of the Data Normalization Module. For a listing of other stages, see [Data Normalization Module](#) on page 267.

### 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 (πσ) 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 (πσ) 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, Transliterater 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 ( $\varsigma\alpha$ ).


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 Transliterater stage can transliterate any string field. You can specify which fields to transliterate in the Transliterater stage options.										
TransliteraterID	<p>Overrides the default transliteration specified in the Transliterater 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><b>Arabic-Latin</b></td> <td>From Arabic to Latin.</td> </tr> <tr> <td><b>Greek-Latin</b></td> <td>From Greek to Latin.</td> </tr> <tr> <td><b>Latin-Hangul</b></td> <td>From Latin to Hangul.</td> </tr> <tr> <td><b>Latin-Katakana</b></td> <td>From Latin to Katakana.</td> </tr> <tr> <td><b>Fullwidth-Halfwidth</b></td> <td>From full width to half width.</td> </tr> </table>	<b>Arabic-Latin</b>	From Arabic to Latin.	<b>Greek-Latin</b>	From Greek to Latin.	<b>Latin-Hangul</b>	From Latin to Hangul.	<b>Latin-Katakana</b>	From Latin to Katakana.	<b>Fullwidth-Halfwidth</b>	From full width to half width.
<b>Arabic-Latin</b>	From Arabic to Latin.										
<b>Greek-Latin</b>	From Greek to Latin.										
<b>Latin-Hangul</b>	From Latin to Hangul.										
<b>Latin-Katakana</b>	From Latin to Katakana.										
<b>Fullwidth-Halfwidth</b>	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 <a href="#">Transliterator</a> on page 279.</p> <p><b>Note:</b> The Transliterator stage does not support transliteration between all scripts. The <b>From</b> and <b>To</b> 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 <a href="#">Transliterator</a> on page 279.</p> <p><b>Note:</b> The Transliterator stage does not support transliteration between all scripts. The <b>From</b> and <b>To</b> fields automatically reflect the valid values based on your selection.</p>
Swap button 	<p>Click the swap button to exchange the languages in the <b>From</b> and <b>To</b> 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 Module

## Universal Name Module

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

The Universal Name Module consists of:

- **Name Parser (Deprecated)**—This component breaks down personal names, 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.

**Attention:** The Name Parser stage is deprecated and may not be supported in future releases. Use Open Name Parser for parsing names.

- **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 multi-record 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 pre-define 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><b>SLAVIC</b></td> <td>Bosnia, Poland, Albania.</td> </tr> <tr> <td><b>ARMENIAN</b></td> <td>Armenia.</td> </tr> <tr> <td><b>DEFAULT</b></td> <td>Bulgaria, Cayman Islands, Ireland, U.S., U.K.</td> </tr> <tr> <td><b>FRENCH</b></td> <td>France.</td> </tr> <tr> <td><b>SCANDINAVIAN</b></td> <td>Denmark, Finland, Iceland, Norway, Sweden.</td> </tr> <tr> <td><b>GERMANIC</b></td> <td>Austria, Germany, Luxembourg, Switzerland, The Netherlands.</td> </tr> <tr> <td><b>GREEK</b></td> <td>Greece.</td> </tr> <tr> <td><b>HUNGARIAN</b></td> <td>Hungary.</td> </tr> <tr> <td><b>ITALIAN</b></td> <td>Italy.</td> </tr> <tr> <td><b>PORTUGUESE</b></td> <td>Portugal.</td> </tr> <tr> <td><b>ROMANIA</b></td> <td>Romania.</td> </tr> <tr> <td><b>HISPANIC</b></td> <td>Spain.</td> </tr> <tr> <td><b>ARABIC</b></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 <a href="#">Name Variant Finder</a> on page 306.</p>	<b>SLAVIC</b>	Bosnia, Poland, Albania.	<b>ARMENIAN</b>	Armenia.	<b>DEFAULT</b>	Bulgaria, Cayman Islands, Ireland, U.S., U.K.	<b>FRENCH</b>	France.	<b>SCANDINAVIAN</b>	Denmark, Finland, Iceland, Norway, Sweden.	<b>GERMANIC</b>	Austria, Germany, Luxembourg, Switzerland, The Netherlands.	<b>GREEK</b>	Greece.	<b>HUNGARIAN</b>	Hungary.	<b>ITALIAN</b>	Italy.	<b>PORTUGUESE</b>	Portugal.	<b>ROMANIA</b>	Romania.	<b>HISPANIC</b>	Spain.	<b>ARABIC</b>	Tunisia.
<b>SLAVIC</b>	Bosnia, Poland, Albania.																										
<b>ARMENIAN</b>	Armenia.																										
<b>DEFAULT</b>	Bulgaria, Cayman Islands, Ireland, U.S., U.K.																										
<b>FRENCH</b>	France.																										
<b>SCANDINAVIAN</b>	Denmark, Finland, Iceland, Norway, Sweden.																										
<b>GERMANIC</b>	Austria, Germany, Luxembourg, Switzerland, The Netherlands.																										
<b>GREEK</b>	Greece.																										
<b>HUNGARIAN</b>	Hungary.																										
<b>ITALIAN</b>	Italy.																										
<b>PORTUGUESE</b>	Portugal.																										
<b>ROMANIA</b>	Romania.																										
<b>HISPANIC</b>	Spain.																										
<b>ARABIC</b>	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 <b>Remove</b> .	<p>Click this box to separate names containing more than one individual into multiple records, for example, <i>Bill &amp; Sally Smith</i>.</p> <p>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.</p>
Gender Determination Source Select a match results in the Match Results List and then click <b>Remove</b> .	<p>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.</p> <p><b>Note:</b> 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.</p>
Order	<p>Specifies how the name fields are ordered in your input records. One of the following:</p> <p><b>Natural</b>            The name fields are ordered by Title, First Name, Middle Name, Last Name, and Suffix.</p> <p><b>Reverse</b>            The name fields are ordered by Last Name first.</p> <p><b>Mixed</b>                The name fields are ordered using a combination of natural and reverse.</p>
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 <a href="#">Modifying Name Parser User-Defined Tables</a> on page 289.

### **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\Pitney Bowes\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 <a href="#">Input</a> on page 286.

Column Name	Description / Valid Values
Gender	<p>The gender most commonly associated with this FirstName/Culture combination. One of the following:</p> <p><b>M</b>      The name is a male name.</p> <p><b>F</b>      The name is a female name.</p> <p><b>A</b>      Ambiguous. The name can be either male or female.</p> <p><b>U</b>      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>
</table-data>
```

```

    ]]>
  </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><b>M</b>      The name is a male name.</p> <p><b>F</b>      The name is a female name.</p> <p><b>A</b>      Ambiguous. The name can be either male or female.</p> <p><b>U</b>      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 <a href="#">Input</a> on page 286.

## 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
        ALII|M|DEFAULT
      ]]>
    </deleted-entry-group>
  </deleted-entries>
</table-data>
```

```

        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>
</table-data>

```



```

        <![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>
```

### *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><b>M</b>      The name is a male name.</p> <p><b>F</b>      The name is a female name.</p> <p><b>A</b>      Ambiguous. The name can be either male or female.</p> <p><b>U</b>      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 <a href="#">Input</a> on page 286.

#### 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
      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"> <li><b>M</b>      The name is a male name.</li> <li><b>F</b>      The name is a female name.</li> <li><b>A</b>      Ambiguous. The name can be either male or female.</li> <li><b>U</b>      Unknown. The gender of this name is not known. Unknown is assumed if this field is left blank.</li> </ul>

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
      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: <b>Firm</b> The name is a company name. <b>Personal</b> The name is an individual person's name.
<b>Fields Related to Names of Companies</b>		
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".
FirmName	String	The name of a company. For example, "Pitney Bowes, Inc."

Field Name	Format	Description / Valid Values
FirmPrimary	String	The base part of a company's name. For example, "Pitney Bowes".
FirmSuffix	String	The corporate suffix. For example, "Co." and "Inc."
<b>Fields Related to Names of Individual People</b>		
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> <p><b>A</b>      Ambiguous. The name is both a male and a female name. For example, Pat.</p> <p><b>F</b>      Female. The name is a female name.</p> <p><b>M</b>      Male. The name is a male name.</p> <p><b>U</b>      Unknown. The name could not be found in the gender table.</p>
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.
MaturitySuffix	String	A person's maturity/generational suffix. For example, Jr. or Sr.

Field Name	Format	Description / Valid Values
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.
<b>Fields Related to Conjoined Names</b>		
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><b>A</b> Ambiguous. The name is both a male and a female name. For example, Pat.</p> <p><b>F</b> Female. The name is a female name.</p> <p><b>M</b> Male. The name is a male name.</p> <p><b>U</b> Unknown. The name could not be found in the gender table.</p>



Field Name	Format	Description / Valid Values
PersonalName.2.GenderDeterminationSource	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. &amp; Mrs. John Smith &amp; Adam Jones". One of the following:</p> <p><b>A</b> Ambiguous. The name is both a male and a female name. For example, Pat.</p> <p><b>F</b> Female. The name is a female name.</p> <p><b>M</b> Male. The name is a male name.</p> <p><b>U</b> 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 Universal Name Module, Data Normalization Module, and Advanced Matching Modules 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><b>Note:</b> Gender codes only apply to first names, not last names.</p> <table border="0"> <tr> <td><b>M</b></td> <td>The name is a male name.</td> </tr> <tr> <td><b>F</b></td> <td>The name is a female name.</td> </tr> <tr> <td><b>A</b></td> <td>Ambiguous. The name can be either male or female.</td> </tr> <tr> <td><b>U</b></td> <td>Unknown. The gender of this name is not known.</td> </tr> </table>	<b>M</b>	The name is a male name.	<b>F</b>	The name is a female name.	<b>A</b>	Ambiguous. The name can be either male or female.	<b>U</b>	Unknown. The gender of this name is not known.
<b>M</b>	The name is a male name.								
<b>F</b>	The name is a female name.								
<b>A</b>	Ambiguous. The name can be either male or female.								
<b>U</b>	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><b>Note:</b> 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.

Option	Description
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 <a href="#">Input</a> on page 307.
Ethnicity	Returns name variations only for the culture specified in the record's Ethnicity field. For information about the Ethnicity field, see <a href="#">Input</a> on page 307.
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 <b>Native</b> , you can choose to return Japanese names in Kana by selecting this option. Kana is comprised of hiragana and katakana scripts.  <b>Note:</b> 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 <b>Native</b> , you can choose to return Japanese names in Kanji by selecting this option. Kanji is one of the scripts used in the Japanese language.  <b>Note:</b> 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.  <b>Note:</b> 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:  <b>M</b> The name is a male name. <b>F</b> The name is a female name. <b>A</b> Ambiguous. The name can be either male or female. <b>U</b> 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:  <b>Suspect</b> A suspect record is used as input to a query. <b>Candidate</b> 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 multi-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><b>Null (empty)</b></td> <td>Global culture (default).</td> </tr> <tr> <td><b>de</b></td> <td>German.</td> </tr> <tr> <td><b>es</b></td> <td>Spanish.</td> </tr> <tr> <td><b>ja</b></td> <td>Japanese.</td> </tr> </table> <p><b>Note:</b> If you added your own domain using the Open Parser Domain Editor, the cultures and culture codes for that domain are also valid.</p>	<b>Null (empty)</b>	Global culture (default).	<b>de</b>	German.	<b>es</b>	Spanish.	<b>ja</b>	Japanese.
<b>Null (empty)</b>	Global culture (default).								
<b>de</b>	German.								
<b>es</b>	Spanish.								
<b>ja</b>	Japanese.								

Field Name	Description
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><b>Natural</b>     The name fields are ordered by Title, First Name, Middle Name, Last Name, and Suffix.</p> <p><b>Reverse</b>     The name fields are ordered by Last Name first.</p> <p><b>Both</b>     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 &amp; 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.

Option Name	Description
Output results as list	Specifies whether to return the parsed name elements in a list form.
Shortcut threshold	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.  The default is 100.

### *Cultures Options*

The following table lists the options that control name cultures.

**Table 48: Open Name Parser Cultures Options**

Option Name	Description
Cultures	Specifies which culture(s) you want to include in the parsing grammar. Global Culture is the default selection.  <b>Note:</b> 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.  Click the <b>Up</b> and <b>Down</b> buttons to set the order in which you want the cultures to run.

### *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 <b>Up</b> and <b>Down</b> 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 <b>Shortcut threshold</b> 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><b>Note:</b> 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 <b>Output results as list</b> box under Parsing Options.
<b>Fields Related to Names of Companies</b>		
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, "Pitney Bowes".
FirmSuffix	String	The corporate suffix. For example, "Co." and "Inc."
IsFirm	String	Indicates that the name is a firm rather than an individual.
<b>Fields Related to Names of Individual People</b>		
Conjunction	String	Indicates that the name contains a conjunction such as "and", "or", or "&".
CultureCode	String	The culture codes contained in the input data.

Field Name	Format	Description
CultureCodeUsedToParse	String	<p>Identifies the culture-specific grammar that was used to parse the data.</p> <p><b>Null (empty)</b> Global culture (default).</p> <p><b>de</b> German.</p> <p><b>es</b> Spanish.</p> <p><b>ja</b> Japanese.</p> <p><b>Note:</b> 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.</p>
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.
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.

Field Name	Format	Description
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.
<b>Fields Related to Conjoined Names</b>		
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.
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.

Field Name	Format	Description
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."
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".

# 9 - ISO Country Codes and Module Support

In this section

---

ISO Country Codes and Module Support

320

# 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 Enterprise Geocoding Module 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 3116-1 Alpha-2	ISO 3116-1 Alpha-3	Supported Modules
Afghanistan	AF	AFG	Universal Addressing Module
Aland Islands	AX	ALA	Universal Addressing Module
Albania	AL or SQ (Routing)	ALB	Universal Addressing Module Enterprise Geocoding Module Enterprise Routing Module
Algeria	DZ	DZA	Enterprise Geocoding Module (Africa) Universal Addressing Module
American Samoa	AS	ASM	Universal Addressing Module
Andorra	AD	AND	Enterprise Geocoding Module. (Andorra is covered by the Spain geocoder) Universal Addressing Module



ISO Country Name	ISO 3116-1 Alpha-2	ISO 3116-1 Alpha-3	Supported Modules
Angola	AO	AGO	Enterprise Geocoding Module (Africa) Universal Addressing Module
Anguilla	AI	AIA	Universal Addressing Module
Antarctica	AQ	ATA	Universal Addressing Module
Antigua And Barbuda	AG	ATG	Universal Addressing Module
Argentina	AR	ARG	Enterprise Geocoding Module Universal Addressing Module Enterprise Routing Module
Armenia	AM	ARM	Universal Addressing Module
Aruba	AW	ABW	Enterprise Geocoding Module (Latin America) Universal Addressing Module
Australia	AU	AUS	Enterprise Geocoding Module Enterprise Routing Module Universal Addressing Module
Austria	AT	AUT	Enterprise Geocoding Module Enterprise Routing Module Universal Addressing Module
Azerbaijan	AZ	AZE	Universal Addressing Module

ISO Country Name	ISO 3116-1 Alpha-2	ISO 3116-1 Alpha-3	Supported Modules
Bahamas	BS	BHS	Enterprise Geocoding Module Universal Addressing Module Enterprise Routing Module
Bahrain	BH	BHR	Enterprise Geocoding Module (Middle East) Universal Addressing Module
Bangladesh	BD	BGD	Universal Addressing Module
Barbados	BB	BRB	Enterprise Geocoding Module (Latin America) Universal Addressing Module
Belarus	BY	BLR	Universal Addressing Module Enterprise Routing Module
Belgium	BE	BEL	Enterprise Geocoding Module Enterprise Routing Module Universal Addressing Module
Belize	BZ	BLZ	Enterprise Geocoding Module (Latin America) Universal Addressing Module
Benin	BJ	BEN	Enterprise Geocoding Module (Africa) Universal Addressing Module
Bermuda	BM	BMU	Universal Addressing Module Enterprise Routing Module
Bhutan	BT	BTN	Universal Addressing Module

ISO Country Name	ISO 3116-1 Alpha-2	ISO 3116-1 Alpha-3	Supported Modules
Bolivia	BO	BOL	Enterprise Geocoding Module (Latin America) Universal Addressing Module
Bonaire, Saint Eustatius And Saba	BQ	BES	Universal Addressing Module
Bosnia And Herzegovina	BA	BIH	Enterprise Geocoding Module Universal Addressing Module Enterprise Routing Module Enterprise Geocoding Module
Botswana	BW	BWA	Enterprise Geocoding Module (Africa) Universal Addressing Module
Bouvet Island	BV	BVT	Universal Addressing Module
Brazil	BR	BRA	Enterprise Geocoding Module Enterprise Routing Module Universal Addressing Module
British Indian Ocean Territory	IO	IOT	Universal Addressing Module
Brunei Darussalam	BN	BRN	Enterprise Geocoding Module Universal Addressing Module
Bulgaria	BG	BGR	Enterprise Geocoding Module Universal Addressing Module
Burkina Faso	BF	BFA	Enterprise Geocoding Module (Africa) Universal Addressing Module

ISO Country Name	ISO 3116-1 Alpha-2	ISO 3116-1 Alpha-3	Supported Modules
Burundi	BI	BDI	Enterprise Geocoding Module (Africa) Universal Addressing Module
Cambodia	KH	KHM	Universal Addressing Module
Cameroon	CM	CMR	Enterprise Geocoding Module (Africa) Universal Addressing Module
Canada	CA	CAN	Enterprise Geocoding Module Enterprise Routing Module Universal Addressing Module
Cape Verde	CV	CPV	Universal Addressing Module
Cayman Islands	KY	CYM	Universal Addressing Module
Central African Republic	CF	CAF	Universal Addressing Module
Chad	TD	TCD	Universal Addressing Module
Chile	CL	CHL	Enterprise Geocoding Module Universal Addressing Module Enterprise Routing Module
China	CN or zh_CN (Routing)	CHN	Enterprise Geocoding Module Universal Addressing Module Enterprise Routing Module
Christmas Island	CX	CXR	Universal Addressing Module

ISO Country Name	ISO 3116-1 Alpha-2	ISO 3116-1 Alpha-3	Supported Modules
Cocos (Keeling) Islands	CC	CCK	Universal Addressing Module
Colombia	CO	COL	Enterprise Geocoding Module Universal Addressing Module
Comoros	KM	COM	Universal Addressing Module
Congo, Republic Of The	CG	COG	Enterprise Geocoding Module (Africa) Universal Addressing Module
Congo, The Democratic Republic Of The	CD	COD	Enterprise Geocoding Module (Africa) Universal Addressing Module Enterprise Routing Module
Cook Islands	CK	COK	Universal Addressing Module
Costa Rica	CR	CRI	Enterprise Geocoding Module (Latin America) Universal Addressing Module
Côte d'Ivoire	CI	CIV	Universal Addressing Module
Croatia	HR	HRV	Enterprise Geocoding Module Enterprise Routing Module Universal Addressing Module
Cuba	CU	CUB	Enterprise Geocoding Module (Latin America) Enterprise Routing Module Universal Addressing Module
Curacao	CW	CUW	Universal Addressing Module

ISO Country Name	ISO 3116-1 Alpha-2	ISO 3116-1 Alpha-3	Supported Modules
Cyprus	CY	CYP	Enterprise Geocoding Module Universal Addressing Module
Czech Republic	CZ or CS (Routing)	CZE	Enterprise Geocoding Module Universal Addressing Module Enterprise Routing Module
Denmark	DK	DNK	Enterprise Geocoding Module Enterprise Routing Module Universal Addressing Module
Djibouti	DJ	DJI	Universal Addressing Module
Dominica	DM	DMA	Universal Addressing Module
Dominican Republic	DO	DOM	Enterprise Geocoding Module (Latin America) Universal Addressing Module
Ecuador	EC	ECU	Enterprise Geocoding Module (Latin America) Universal Addressing Module
Egypt	EG	EGY	Enterprise Geocoding Module (Middle East) Universal Addressing Module
El Salvador	SV	SLV	Enterprise Geocoding Module (Latin America) Universal Addressing Module
Equatorial Guinea	GQ	GNQ	Universal Addressing Module
Eritrea	ER	ERI	Universal Addressing Module

ISO Country Name	ISO 3116-1 Alpha-2	ISO 3116-1 Alpha-3	Supported Modules
Estonia	EE	EST	Enterprise Geocoding Module Enterprise Routing Module Universal Addressing Module
Ethiopia	ET	ETH	Universal Addressing Module
Falkland Islands (Malvinas)	FK	FLK	Universal Addressing Module
Faroe Islands	FO	FRO	Universal Addressing Module
Fiji	FJ	FJI	Universal Addressing Module
Finland	FI	FIN	Enterprise Geocoding Module Enterprise Routing Module Universal Addressing Module
France	FR	FRA	Enterprise Geocoding Module Enterprise Routing Module Universal Addressing Module
French Guiana	GF	GUF	Enterprise Geocoding Module ( <i>French Guiana is covered by the France geocoder.</i> ) Universal Addressing Module
French Polynesia	PF	PYF	Universal Addressing Module
French Southern Territories	TF	ATF	Universal Addressing Module
Gabon	GA	GAB	Enterprise Geocoding Module (Africa) Universal Addressing Module

ISO Country Name	ISO 3116-1 Alpha-2	ISO 3116-1 Alpha-3	Supported Modules
Gambia	GM	GMB	Universal Addressing Module
Georgia	GE	GEO	Universal Addressing Module
Germany	DE	DEU	Enterprise Geocoding Module Enterprise Routing Module Universal Addressing Module
Ghana	GH	GHA	Enterprise Geocoding Module (Africa) Universal Addressing Module Enterprise Routing Module
Gibraltar	GI	GIB	Enterprise Geocoding Module ( <i>Gibraltar is covered by the Spain geocoder.</i> ) Universal Addressing Module
Greece	GR	GRC	Enterprise Geocoding Module Universal Addressing Module
Greenland	GL	GRL	Universal Addressing Module
Grenada	GD	GRD	Universal Addressing Module
Guadeloupe	GP	GLP	Enterprise Geocoding Module ( <i>Guadeloupe is covered by the France geocoder.</i> ) Universal Addressing Module
Guam	GU	GUM	Universal Addressing Module



ISO Country Name	ISO 3116-1 Alpha-2	ISO 3116-1 Alpha-3	Supported Modules
Guatemala	GT	GTM	Enterprise Geocoding Module (Latin America) Universal Addressing Module
Guernsey	GG	GGY	Universal Addressing Module
Guinea	GN	GIN	Universal Addressing Module
Guinea-Bissau	GW	GNB	Universal Addressing Module
Guyana	GY	GUY	Enterprise Geocoding Module (Latin America) Universal Addressing Module
Haiti	HT	HTI	Universal Addressing Module
Heard Island and McDonald Islands	HM	HMD	Universal Addressing Module
Holy See (Vatican City State)	VA	VAT	Enterprise Geocoding Module ( <i>The Vatican is covered by the Italy geocoder.</i> ) Universal Addressing Module
Honduras	HN	HND	Enterprise Geocoding Module (Latin America) Universal Addressing Module
Hong Kong	HK	HKG	Enterprise Geocoding Module Universal Addressing Module
Hungary	HU	HUN	Enterprise Geocoding Module Enterprise Routing Module Universal Addressing Module

ISO Country Name	ISO 3116-1 Alpha-2	ISO 3116-1 Alpha-3	Supported Modules
Iceland	IS	ISL	Enterprise Geocoding Module Universal Addressing Module
India	IN	IND	Enterprise Geocoding Module Universal Addressing Module
Indonesia	ID	IDN	Enterprise Geocoding Module Universal Addressing Module
Iran, Islamic Republic Of	IR	IRN	Universal Addressing Module
Iraq	IQ	IRQ	Enterprise Geocoding Module (Middle East) Universal Addressing Module
Ireland	IE	IRL	Enterprise Geocoding Module Enterprise Routing Module Universal Addressing Module
Isle Of Man	IM	IMN	Universal Addressing Module
Israel	IL	ISR	Universal Addressing Module Enterprise Routing Module
Italy	IT	ITA	Enterprise Geocoding Module Enterprise Routing Module Universal Addressing Module
Jamaica	JM	JAM	Enterprise Geocoding Module (Latin America) Universal Addressing Module

ISO Country Name	ISO 3116-1 Alpha-2	ISO 3116-1 Alpha-3	Supported Modules
Japan	JP	JPN	Enterprise Geocoding Module Enterprise Routing Module Universal Addressing Module
Jersey	JE	JEY	Universal Addressing Module
Jordan	JO	JOR	Universal Addressing Module Enterprise Geocoding Module (Middle East) Enterprise Routing Module
Kazakhstan	KZ	KAZ	Universal Addressing Module
Kenya	KE	KEN	Enterprise Geocoding Module (Africa) Universal Addressing Module Enterprise Routing Module
Kiribati	KI	KIR	Universal Addressing Module
Korea, Democratic People's Republic Of	KP	PRK	Universal Addressing Module
Korea, Republic Of	KR	KOR	Enterprise Geocoding Module Universal Addressing Module
Kosovo	Xk	XKX	Enterprise Geocoding Module Universal Addressing Module
Kuwait	KW	KWT	Enterprise Geocoding Module (Middle East) Universal Addressing Module

ISO Country Name	ISO 3116-1 Alpha-2	ISO 3116-1 Alpha-3	Supported Modules
Kyrgyzstan	KG	KGZ	Universal Addressing Module
Lao People's Democratic Republic	LA	LAO	Universal Addressing Module
Latvia	LV	LVA	Enterprise Geocoding Module Enterprise Routing Module Universal Addressing Module
Lebanon	LB	LBN	Enterprise Geocoding Module (Middle East) Universal Addressing Module
Lesotho	LS	LSO	Enterprise Geocoding Module (Africa) Universal Addressing Module Enterprise Routing Module
Liberia	LR	LBR	Universal Addressing Module
Libyan Arab Jamahiriya	LY	LYB	Universal Addressing Module
Liechtenstein	LI	LIE	Enterprise Geocoding Module ( <i>Liechtenstein is covered by the Switzerland geocoder.</i> ) Enterprise Routing Module Universal Addressing Module
Lithuania	LT	LTU	Enterprise Geocoding Module Enterprise Routing Module Universal Addressing Module
Luxembourg	LU	LUX	Enterprise Geocoding Module ( <i>Luxembourg is covered by the Belgium geocoder.</i> ) Enterprise Routing Module Universal Addressing Module

ISO Country Name	ISO 3116-1 Alpha-2	ISO 3116-1 Alpha-3	Supported Modules
Macao	MO	MAC	Enterprise Geocoding Module Universal Addressing Module
Macedonia, Former Yugoslav Republic Of	MK	MKD	Enterprise Geocoding Module Universal Addressing Module
Madagascar	MG	MDG	Universal Addressing Module
Malawi	MW	MWI	Enterprise Geocoding Module (Africa) Universal Addressing Module
Malaysia	MY	MYS	Enterprise Geocoding Module Enterprise Routing Module Universal Addressing Module
Maldives	MV	MDV	Universal Addressing Module
Mali	ML	MLI	Enterprise Geocoding Module (Africa) Universal Addressing Module
Malta	ML	MLT	Enterprise Geocoding Module Universal Addressing Module
Marshall Islands	MH	MHL	Universal Addressing Module
Martinique	MQ	MTQ	Enterprise Geocoding Module ( <i>Martinique is covered by the France geocoder.</i> ) Universal Addressing Module
Mauritania	MR	MRT	Enterprise Geocoding Module (Africa) Universal Addressing Module

ISO Country Name	ISO 3116-1 Alpha-2	ISO 3116-1 Alpha-3	Supported Modules
Mauritius	MU	MUS	Enterprise Geocoding Module (Africa) Universal Addressing Module
Mayotte	YT	MYT	Enterprise Geocoding Module ( <i>Mayotte is covered by the France geocoder.</i> ) Universal Addressing Module
Mexico	MX	MEX	Enterprise Geocoding Module Universal Addressing Module
Micronesia, Federated States Of	FM	FSM	Universal Addressing Module
Moldova, Republic Of	MD	MDA	Universal Addressing Module Enterprise Routing Module
Monaco	MC	MCO	Enterprise Geocoding Module ( <i>Monaco is covered by the France geocoder.</i> ) Universal Addressing Module
Mongolia	MN	MNG	Universal Addressing Module
Montenegro	ME	MNE	Enterprise Geocoding Module Universal Addressing Module
Montserrat	MS	MSR	Universal Addressing Module
Morocco	MA	MAR	Enterprise Geocoding Module (Africa) Universal Addressing Module

ISO Country Name	ISO 3116-1 Alpha-2	ISO 3116-1 Alpha-3	Supported Modules
Mozambique	MZ	MOZ	Enterprise Geocoding Module (Africa) Universal Addressing Module Enterprise Routing Module
Myanmar	MM	MMR	Universal Addressing Module
Namibia	NA	NAM	Enterprise Geocoding Module (Africa) Universal Addressing Module
Nauru	NR	NRU	Universal Addressing Module
Nepal	NP	NPL	Universal Addressing Module
Netherlands	NL	NLD	Enterprise Geocoding Module Enterprise Routing Module Universal Addressing Module
New Caledonia	NC	NCL	Universal Addressing Module
New Zealand	NZ	NZL	Enterprise Geocoding Module Enterprise Routing Module Universal Addressing Module
Nicaragua	NI	NIC	Enterprise Geocoding Module (Latin America) Universal Addressing Module
Niger	NE	NER	Enterprise Geocoding Module (Africa) Universal Addressing Module

ISO Country Name	ISO 3116-1 Alpha-2	ISO 3116-1 Alpha-3	Supported Modules
Nigeria	NG	NGA	Enterprise Geocoding Module (Africa) Universal Addressing Module
Niue	NU	NIU	Universal Addressing Module
Norfolk Island	NF	NFK	Universal Addressing Module
Northern Mariana Islands	MP	MNP	Universal Addressing Module
Norway	NO	NOR	Enterprise Geocoding Module Enterprise Routing Module Universal Addressing Module
Oman	OM	OMN	Enterprise Geocoding Module (Middle East) Universal Addressing Module
Pakistan	PK	PAK	Universal Addressing Module
Palau	PW	PLW	Universal Addressing Module
Palestinian Territory, Occupied	PS	PSE	Universal Addressing Module
Panama	PA	PAN	Enterprise Geocoding Module (Latin America) Universal Addressing Module
Papua New Guinea	PG	PNG	Universal Addressing Module



ISO Country Name	ISO 3116-1 Alpha-2	ISO 3116-1 Alpha-3	Supported Modules
Paraguay	PY	PRY	Enterprise Geocoding Module (Latin America) Universal Addressing Module
Peru	PE	PER	Enterprise Geocoding Module (Latin America) Universal Addressing Module
Philippines	PH	PHL	Enterprise Geocoding Module Universal Addressing Module Enterprise Routing Module
Pitcairn	PN	PCN	Universal Addressing Module
Poland	PL	POL	Enterprise Geocoding Module Enterprise Routing Module Universal Addressing Module
Portugal	PT	PRT	Enterprise Geocoding Module Enterprise Routing Module Universal Addressing Module
Puerto Rico	PR	PRI	Universal Addressing Module
Qatar	QA	QAT	Enterprise Geocoding Module (Middle East) Universal Addressing Module
Reunion	RE	REU	Enterprise Geocoding Module ( <i>Reunion is covered by the France geocoder.</i> ) Universal Addressing Module
Romania	RO	ROU	Enterprise Geocoding Module Enterprise Routing Module Universal Addressing Module

ISO Country Name	ISO 3116-1 Alpha-2	ISO 3116-1 Alpha-3	Supported Modules
Russian Federation	RU	RUS	Enterprise Geocoding Module Universal Addressing Module
Rwanda	RW	RWA	Enterprise Geocoding Module (Africa) Universal Addressing Module
Saint Barthelemy	BL	BLM	Universal Addressing Module
Saint Helena, Ascension and Tristan Da Cunha	SH	SHE	Universal Addressing Module
Saint Kitts and Nevis	KN	KNA	Enterprise Geocoding Module (Latin America) Universal Addressing Module
Saint Lucia	LC	LCA	Universal Addressing Module
Saint Martin (French Part)	MF	MAF	Universal Addressing Module
Saint Pierre and Miquelon	PM	SPM	Universal Addressing Module
Saint Vincent and the Grenadines	VC	VCT	Universal Addressing Module
Samoa	WS	WSM	Universal Addressing Module
San Marino	SM	SMR	Enterprise Geocoding Module ( <i>San Marino is covered by the Italy geocoder.</i> ) Universal Addressing Module
Sao Tome and Principe	ST	STP	Universal Addressing Module

ISO Country Name	ISO 3116-1 Alpha-2	ISO 3116-1 Alpha-3	Supported Modules
Saudi Arabia	SA	SAU	Enterprise Geocoding Module (Middle East) Universal Addressing Module
Senegal	SN	SEN	Enterprise Geocoding Module (Africa) Universal Addressing Module
Serbia	RS	SRB	Enterprise Geocoding Module Universal Addressing Module
Seychelles	SC	SYC	Universal Addressing Module
Sierra Leone	SL	SLE	Universal Addressing Module
Singapore	SG	SGP	Enterprise Geocoding Module Enterprise Routing Module Universal Addressing Module
Sint Maarten (Dutch Part)	SX	SXM	Universal Addressing Module
Slovakia	SK	SVK	Enterprise Geocoding Module Enterprise Routing Module Universal Addressing Module
Slovenia	SI	SVN	Enterprise Geocoding Module Enterprise Routing Module Universal Addressing Module
Solomon Islands	SB	SLB	Universal Addressing Module
Somalia	SO	SOM	Universal Addressing Module

ISO Country Name	ISO 3116-1 Alpha-2	ISO 3116-1 Alpha-3	Supported Modules
South Africa	ZA	ZAF	Enterprise Geocoding Module Universal Addressing Module
South Georgia And The South Sandwich Islands	GS	SGS	Enterprise Geocoding Module Universal Addressing Module
South Sudan	SS	SSD	Universal Addressing Module
Spain	ES	ESP	Enterprise Geocoding Module Enterprise Routing Module Universal Addressing Module
Sri Lanka	LK	LKA	Universal Addressing Module
Sudan	SD	SDN	Universal Addressing Module
Suriname	SR	SUR	Enterprise Geocoding Module (Latin America) Universal Addressing Module
Svalbard And Jan Mayen	SJ	SJM	Universal Addressing Module
Swaziland	SZ	SWZ	Enterprise Geocoding Module (Africa) Universal Addressing Module
Sweden	SE	SWE	Enterprise Geocoding Module Enterprise Routing Module Universal Addressing Module
Switzerland	CH	CHE	Enterprise Geocoding Module Enterprise Routing Module Universal Addressing Module

ISO Country Name	ISO 3116-1 Alpha-2	ISO 3116-1 Alpha-3	Supported Modules
Syrian Arab Republic	SY	SYR	Universal Addressing Module
Taiwan, Province of China	TW or zh_TW (Routing)	TWN	Universal Addressing Module Enterprise Routing Module
Tajikistan	TJ	TJK	Universal Addressing Module
Tanzania, United Republic Of	TZ	TZA	Enterprise Geocoding Module (Africa) Universal Addressing Module Enterprise Routing Module
Thailand	TH	THA	Enterprise Geocoding Module Enterprise Routing Module Universal Addressing Module
Timor-Leste	TL	TLS	Universal Addressing Module
Togo	TG	TGO	Enterprise Geocoding Module (Africa) Universal Addressing Module
Tokelau	TK	TKL	Universal Addressing Module
Tonga	TO	TON	Universal Addressing Module
Trinidad and Tobago	TT	TTO	Enterprise Geocoding Module (Latin America) Universal Addressing Module
Tunisia	TN	TUN	Enterprise Geocoding Module (Africa) Universal Addressing Module

ISO Country Name	ISO 3116-1 Alpha-2	ISO 3116-1 Alpha-3	Supported Modules
Turkey	TR	TUR	Enterprise Geocoding Module Universal Addressing Module
Turkmenistan	TM	TKM	Universal Addressing Module
Turks And Caicos Islands	TC	TCA	Universal Addressing Module
Tuvalu	TV	TUV	Universal Addressing Module
Uganda	UG	UGA	Enterprise Geocoding Module (Africa) Universal Addressing Module
Ukraine	UA	UKR	Enterprise Geocoding Module Universal Addressing Module
United Arab Emirates	AE	ARE	Enterprise Geocoding Module (Middle East) Universal Addressing Module
United Kingdom	GB	GBR	Enterprise Geocoding Module Enterprise Routing Module Universal Addressing Module
United States	US	USA	Enterprise Geocoding Module Enterprise Routing Module Universal Addressing Module
United States Minor Outlying Islands	UM	UMI	Universal Addressing Module
Uruguay	UY	URY	Enterprise Geocoding Module Universal Addressing Module

ISO Country Name	ISO 3116-1 Alpha-2	ISO 3116-1 Alpha-3	Supported Modules
Uzbekistan	UZ	UZB	Universal Addressing Module
Vanuatu	VU	VUT	Universal Addressing Module
Venezuela, Bolivarian Republic Of	VE	VEN	Enterprise Geocoding Module Universal Addressing Module
Viet Nam	VN	VNM	Enterprise Geocoding Module Universal Addressing Module
Virgin Islands, British	VG	VGB	Universal Addressing Module
Virgin Islands, U.S.	VI	VIR	Universal Addressing Module
Wallis and Futuna	WF	WLF	Universal Addressing Module
Western Sahara	EH	ESH	Universal Addressing Module
Yemen	YE	YEM	Enterprise Geocoding Module (Middle East) Universal Addressing Module
Zambia	ZM	ZMB	Enterprise Geocoding Module (Africa) Universal Addressing Module
Zimbabwe	ZW	ZWE	Enterprise Geocoding Module (Africa) Universal Addressing Module

# Notices



© 2018 Pitney Bowes. All rights reserved. MapInfo and Group 1 Software are trademarks of Pitney Bowes Software Inc. All other marks and trademarks are property of their respective holders.

### *USPS® Notices*

Pitney Bowes Inc. holds a non-exclusive license to publish and sell ZIP + 4® databases on optical and magnetic media. The following trademarks are owned by the United States Postal Service: CASS, CASS Certified, DPV, eLOT, FASTforward, First-Class Mail, Intelligent Mail, LACS<sup>Link</sup>, NCOA<sup>Link</sup>, PAVE, PLANET Code, Postal Service, POSTNET, Post Office, RDI, Suite<sup>Link</sup>, United States Postal Service, Standard Mail, United States Post Office, USPS, ZIP Code, and ZIP + 4. This list is not exhaustive of the trademarks belonging to the Postal Service.

Pitney Bowes Inc. is a non-exclusive licensee of USPS® for NCOA<sup>Link</sup>® processing.

Prices for Pitney Bowes Software's products, options, and services are not established, controlled, or approved by USPS® or United States Government. When utilizing RDI™ data to determine parcel-shipping costs, the business decision on which parcel delivery company to use is not made by the USPS® or United States Government.

### *Data Provider and Related Notices*

Data Products contained on this media and used within Pitney Bowes Software applications are protected by various trademarks and by one or more of the following copyrights:

© Copyright United States Postal Service. All rights reserved.

© 2014 TomTom. All rights reserved. TomTom and the TomTom logo are registered trademarks of TomTom N.V.

© 2016 HERE

Fuente: INEGI (Instituto Nacional de Estadística y Geografía)

Based upon electronic data © National Land Survey Sweden.

© Copyright United States Census Bureau

© Copyright Nova Marketing Group, Inc.

Portions of this program are © Copyright 1993-2007 by Nova Marketing Group Inc. All Rights Reserved

© Copyright Second Decimal, LLC

© Copyright Canada Post Corporation

This CD-ROM contains data from a compilation in which Canada Post Corporation is the copyright owner.

© 2007 Claritas, Inc.

The Geocode Address World data set contains data licensed from the GeoNames Project ([www.geonames.org](http://www.geonames.org)) provided under the Creative Commons Attribution License ("Attribution License") located at <http://creativecommons.org/licenses/by/3.0/legalcode>. Your use of the

GeoNames data (described in the Spectrum™ Technology Platform User Manual) is governed by the terms of the Attribution License, and any conflict between your agreement with Pitney Bowes Software, Inc. and the Attribution License will be resolved in favor of the Attribution License solely as it relates to your use of the GeoNames data.



3001 Summer Street  
Stamford CT 06926-0700  
USA

[www.pitneybowes.com](http://www.pitneybowes.com)