



# Understanding the Sun Match Engine



Sun Microsystems, Inc.  
4150 Network Circle  
Santa Clara, CA 95054  
U.S.A.

Part No: 820-2652-15  
December 2008

Copyright 2008 Sun Microsystems, Inc. 4150 Network Circle, Santa Clara, CA 95054 U.S.A. All rights reserved.

Sun Microsystems, Inc. has intellectual property rights relating to technology embodied in the product that is described in this document. In particular, and without limitation, these intellectual property rights may include one or more U.S. patents or pending patent applications in the U.S. and in other countries.

U.S. Government Rights – Commercial software. Government users are subject to the Sun Microsystems, Inc. standard license agreement and applicable provisions of the FAR and its supplements.

This distribution may include materials developed by third parties.

Parts of the product may be derived from Berkeley BSD systems, licensed from the University of California. UNIX is a registered trademark in the U.S. and other countries, exclusively licensed through X/Open Company, Ltd.

Sun, Sun Microsystems, the Sun logo, the Solaris logo, the Java Coffee Cup logo, docs.sun.com, Java, and Solaris are trademarks or registered trademarks of Sun Microsystems, Inc. or its subsidiaries in the U.S. and other countries. All SPARC trademarks are used under license and are trademarks or registered trademarks of SPARC International, Inc. in the U.S. and other countries. Products bearing SPARC trademarks are based upon an architecture developed by Sun Microsystems, Inc.

The OPEN LOOK and Sun<sup>TM</sup> Graphical User Interface was developed by Sun Microsystems, Inc. for its users and licensees. Sun acknowledges the pioneering efforts of Xerox in researching and developing the concept of visual or graphical user interfaces for the computer industry. Sun holds a non-exclusive license from Xerox to the Xerox Graphical User Interface, which license also covers Sun's licensees who implement OPEN LOOK GUIs and otherwise comply with Sun's written license agreements.

Products covered by and information contained in this publication are controlled by U.S. Export Control laws and may be subject to the export or import laws in other countries. Nuclear, missile, chemical or biological weapons or nuclear maritime end uses or end users, whether direct or indirect, are strictly prohibited. Export or reexport to countries subject to U.S. embargo or to entities identified on U.S. export exclusion lists, including, but not limited to, the denied persons and specially designated nationals lists is strictly prohibited.

DOCUMENTATION IS PROVIDED "AS IS" AND ALL EXPRESS OR IMPLIED CONDITIONS, REPRESENTATIONS AND WARRANTIES, INCLUDING ANY IMPLIED WARRANTY OF MERCHANTABILITY, FITNESS FOR A PARTICULAR PURPOSE OR NON-INFRINGEMENT, ARE DISCLAIMED, EXCEPT TO THE EXTENT THAT SUCH DISCLAIMERS ARE HELD TO BE LEGALLY INVALID.

# Contents

---

<b>Understanding the Sun Match Engine</b> .....	5
Related Topics .....	5
About the Sun Match Engine .....	6
Sun Match Engine Overview .....	6
About the Sun Match Engine Matching Algorithm .....	6
Sun Match Engine Standardization and Matching Process .....	7
Sun Match Engine Data Types .....	7
How the Sun Match Engine Works .....	8
Sun Match Engine Matching Weight Formulation .....	9
Sun Match Engine Standardization Configuration .....	10
Sun Match Engine Standardization File Types .....	10
Sun Match Engine Internationalization .....	11
Sun Match Engine Matching Configuration .....	12
The Sun Match Engine Match Configuration File .....	12
The Match Constants File .....	18
Sun Match Engine and the Sun Match Engine .....	18
Master Index Components and the Sun Match Engine .....	19
Configuring the Master Index Matching Service (Repository) .....	27
Sun Match Engine Person Data Type Configuration .....	32
Sun Match Engine Person Matching Overview .....	32
Sun Match Engine Match Configuration for Person Data .....	34
Sun Match Engine Person Data Standardization Files .....	34
Configuring the Sun Match Engine Standardization Files for Person Data .....	43
Configuring the Master Index Matching Service for Person Data (Repository) .....	43
Sun Match Engine Address Data Type Configuration .....	48
Sun Match Engine Address Matching Overview .....	49
Match Configuration for Address Data (Repository) .....	51
Sun Match Engine Standardization Configuration for Address Data .....	51

Modifying Sun Match Engine Address Data Configuration Files .....	62
Configuring the Matching Service for Address Data (Repository) .....	63
Sun Match Engine Business Names Data Type Configuration .....	67
Sun Match Engine Business Name Matching Overview .....	67
Sun Match Engine Match Configuration for Business Names .....	69
Sun Match Engine Standardization Configuration for Business Names .....	69
Modifying Sun Match Engine Business Name Configuration Files .....	83
Configuring the Matching Service for Business Names (Repository) .....	83
Fine-Tuning Weights and Thresholds for Sun Match Engine (Repository) .....	87
Data Analysis Overview .....	87
Customizing the Match Configuration and Thresholds .....	87
Match Configuration Comparison Functions for Sun Match Engine (Repository) .....	93
Sun Match Engine Comparison Functions .....	93
Sun Match Engine Comparison Function Options .....	103

# Understanding the Sun Match Engine

---

The topics listed here provide conceptual information about the Sun Match Engine and how it standardizes and matches data in a master index application.

Note that Java CAPS includes two versions of Sun Match Engine. Sun Match Engine (Repository) is installed in the Java CAPS repository and provides all the functionality of previous versions in the new Java CAPS environment. Sun Match Engine is a service-enabled version of the master index that is installed directly into NetBeans. It includes all of the features of Sun Match Engine (Repository) plus several new features, like data analysis, data cleansing, data loading, and an improved Data Manager GUI. Both products are components of the Sun Master Data Management (MDM) Suite. This document relates to Sun Match Engine (Repository) only.

- [“Sun Match Engine Overview” on page 6](#)
- [“Sun Match Engine Standardization Configuration” on page 10](#)
- [“Sun Match Engine Matching Configuration” on page 12](#)
- [“Sun Match Engine and the Sun Match Engine” on page 18](#)
- [“Sun Match Engine Person Data Type Configuration” on page 32](#)
- [“Sun Match Engine Address Data Type Configuration” on page 48](#)
- [“Sun Match Engine Business Names Data Type Configuration” on page 67](#)
- [“Fine-Tuning Weights and Thresholds for Sun Match Engine \(Repository\)” on page 87](#)
- [“Match Configuration Comparison Functions for Sun Match Engine \(Repository\)” on page 93](#)

## Related Topics

Several topics provide information and instructions for implementing and using a Repository-based master index application. For a complete list of topics related to working with Sun Match Engine (Repository), see “Related Topics” in *Developing Sun Master Indexes (Repository)*.

## About the Sun Match Engine

The Sun Match Engine includes a standardization engine and a match engine, providing data parsing, data standardization, phonetic encoding, and record matching capabilities for master index applications. Before records can be compared to evaluate the possibility of a match, the data contained in those records must be standardized and in certain cases phonetically encoded or parsed. Once the data is conditioned, the match engine determines a match weight for each field defined for matching. The match weight is based on your configuration of the match engine and the fields on which matching is performed. The composite weight (the sum of weights generated for all match fields in the records) indicates how closely two records match.

## Sun Match Engine Overview

The Sun Match Engine is the standard match engine designed to work with the master index applications created by Sun Match Engine (Repository). It is highly configurable in the Sun Match Engine environment and can be used to match on various types of data.

The Sun Match Engine provides data parsing, data standardization, phonetic encoding, and record matching capabilities for master index applications created by Sun Match Engine. Before records can be compared to evaluate the possibility of a match, the data contained in those records must be standardized and in certain cases phonetically encoded or parsed. Once the data is conditioned, the match engine determines a match weight for each field defined for matching. The match weight is based on your configuration of the match engine and the fields on which matching is performed. The composite weight (the sum of weights generated for all match fields in the records) indicates how closely two records match.

The following topics provide information about the configurable components of the match engine and how the Sun Match Engine standardizes and matches data.

- [“About the Sun Match Engine Matching Algorithm”](#) on page 6
- [“Sun Match Engine Standardization and Matching Process”](#) on page 7
- [“Sun Match Engine Data Types”](#) on page 7
- [“How the Sun Match Engine Works”](#) on page 8
- [“Sun Match Engine Matching Weight Formulation”](#) on page 9

## About the Sun Match Engine Matching Algorithm

The Sun Match Engine compares records containing similar data types by calculating how closely the records match. The resulting comparison weight is either a positive or negative numeric value that represents the degree to which the two sets of data are similar. The match engine relies on probabilistic algorithms to compare data of a given type using a comparison function specific to the type of data being compared. The comparison functions for each matching field are defined in a match configuration file that you can customize for the type of

data you are indexing. The formula used to determine the matching weight is based on either matching and unmatching probabilities or on agreement and disagreement weight ranges.

The Sun Match Engine is also designed to standardize free-form text fields, such as street address fields or business names. This allows the match engine to generate a more accurate weight for free-form data.

## Sun Match Engine Standardization and Matching Process

The Sun Match Engine matching algorithm uses a proven methodology to process and weight records in the master index database. By providing both standardization and matching capabilities, the match engine allows you to condition data prior to matching. You can also use these capabilities to review legacy data prior to loading it into the database. This review helps you determine data anomalies, invalid or default values, and missing fields.

Both matching and standardization occur when two records are analyzed for the probability of a match. Before matching, certain fields are normalized, parsed, or converted into their phonetic values if necessary. The match fields are then analyzed and weighted according to the rules defined in a match configuration file. The weights for each field are combined to determine the overall matching weight for the two records. After the match engine has performed these steps, survivorship is determined by the master index application, based on how the overall matching weight compares to the duplicate and match thresholds of the master index application. These thresholds are configured for the Manager Service in the Threshold file.

## Sun Match Engine Data Types

You can standardize and match on different types of data with the Sun Match Engine. In its default implementation with a master index application, the match engine supports data standardization and matching on the three primary types of data listed below.

- Person Information (described in [“Sun Match Engine and the Sun Match Engine” on page 18](#))
- Street Addresses (described in [“Sun Match Engine Person Data Type Configuration” on page 32](#))
- Business Names (described in [“Sun Match Engine Address Data Type Configuration” on page 48](#))

In addition, the Sun Match Engine provides comparison functions for matching on various types of fields contained within the primary data types, such as numbers, dates, Social Security Numbers, single characters, and so on.

When processing person information, the match engine assumes that each match field is stored in a separate field. For street address and business name processing, the match engine can parse free-form text fields for searching and matching. Each data type requires specific customization to the Match Field file in the master index project.

## How the Sun Match Engine Works

The Sun Match Engine compares two records and returns a match weight indicating the likelihood of a match between the two records. The three primary components of the Sun Match Engine are the configuration files, the standardization engine, and the match engine.

- **Configuration Files** - The Sun Match Engine includes several sets of files that define standardization and matching logic for all supported data types. One set of standardization files is common to all national domains, and one additional set is provided for the following national domains: Australia, France, Great Britain, and the United States. You can customize these files to adapt the standardization and matching logic to your specific needs. The matching configuration file defines the configuration of the matching comparator functions.
- **Standardization Engine** - Standardization involves converting nonstandard data into a standardized form for more accurate and efficient processing. Standardization consists of any one or more of the following actions:
  - **Parsing** - Separating a free-text field into its individual components, such as street address information or a business name.
  - **Normalization** - Changing the value of a field to a standard version, such as changing a nickname to a common name.
  - **Phonetic Encoding** - Changing the value of a field to its phonetic version. The field to be converted can be the original field, a parsed field, a normalized field, or a parsed and normalized field.

Using the person data type, for example, first names such as “Bill” and “Will” are normalized to “William”, which is then phonetically converted. Using the street address data type, street addresses are parsed into their component parts, such as house numbers, street names, and so on. The street name is then phonetically converted. Standardization logic is defined in the standardization engine configuration files and in the StandardizationConfig section of the Match Field file, and is performed prior to assigning match weights.

- **Match Engine** - Matching involves comparing two standardized records and returning a weight that indicates the likelihood of a match between the two records. A higher weight indicates a greater likelihood of a match. Matching criteria and logic are defined in the match engine configuration file. The data fields that are sent to the Sun Match Engine for matching, known as the *match string*, are defined in the MatchingConfig section of the Match Field file. The match engine configuration files define how the match string is standardized and which matching rules to use to process each match field.



## Sun Match Engine Matching Weight Formulation

The Sun Match Engine determines the matching weight between two records by comparing the match string fields between the two records using the rules defined in the match configuration file and taking into account the matching logic specified for each field. The Sun Match Engine can use either matching (m) and unmatching (u) conditional probabilities or agreement and disagreement weight ranges to fine-tune the match process. It uses the underlying algorithm to arrive at a match weight for each match string field. The weight generated for each field in the match string indicates the level of match between each field. The weights assigned to each field are then summed together for a total, composite matching weight between the two records. Agreement and disagreement weight ranges or m-probabilities and u-probabilities are defined in the match configuration file.

The following topics describe probabilities and weights.

- [“Matching and Unmatching Probabilities” on page 9](#)
- [“Agreement and Disagreement Weight Ranges” on page 9](#)

### Matching and Unmatching Probabilities

When matching and unmatching conditional probabilities are used, the match engine uses a logarithmic formula to determine agreement and disagreement weights between fields. The m-probabilities and u-probabilities you specify determine the maximum agreement weight and minimum disagreement weight for each field, and so define the agreement and disagreement weight ranges for each field and for the entire record. These probabilities allow you to specify which fields provide the most reliable matching information and which provide the least. For example, in person matching, the gender field is not as reliable as the SSN field for determining a match since a person’s SSN is more specific. Therefore, the SSN field should have a higher m-probability than the gender field. The more reliable the field, the greater the m-probability for that field should be.

If a field matches between two records, an agreement weight, determined by the logarithmic formula using the m-probability and u-probability, is added to the composite match weight for the record. If the fields disagree, a disagreement weight is subtracted from the composite match weight. m-probabilities and u-probabilities are expressed as double values between one and zero (excluding one and zero) and can have up to 16 decimal points.

### Agreement and Disagreement Weight Ranges

Defining agreement and disagreement weight ranges is a more direct way to implement m-probabilities and u-probabilities. Like probabilities, the maximum agreement and minimum disagreement weights you define for each field allow you to define the relative reliability of each field; however, the match weight has a more linear relationship with the numbers you specify. When you use agreement and disagreement weight ranges to determine the match weight, you define a maximum weight for each field when they are in complete agreement and a minimum weight for when they are in complete disagreement. The Sun Match Engine assigns a matching

weight to each field that falls between the agreement and disagreement weights specified for the field. This provides a more convenient and intuitive representation of conditional probabilities.

Using the SSN and gender field example above, the SSN field would be assigned a higher maximum agreement weight and a lower minimum disagreement weight than the gender field because it is more reliable. If you assign a maximum agreement weight of “10” and two SSNs match, the match weight for that field is “10”. If you assign a minimum disagreement weight of “-10” and two SSNs are in complete disagreement, the match weight for that field is “-10”. Agreement and disagreement weights are expressed as double values and can have up to 16 decimal points.

## Sun Match Engine Standardization Configuration

The standardization configuration files define additional logic used by the Sun Match Engine to standardize specific data types. This logic helps define how fields in incoming records are parsed, standardized, and classified for processing. Standardization files include data patterns files, category files, clues files, key type tables, constants files, and reference files.

The standardization configuration files are stored in the master index project and appear as nodes in the Standardization Engine node of the project. Several standardization files are common to all implementations of the Sun Match Engine, but each national domain uses a subset of unique files. The common files are listed directly under the Standardization Engine node of the master index project; the files unique to each national domain are listed in individual sub-folders under the Standardization Engine node.

The standardization configuration files for the Sun Match Engine must follow certain rules for formatting and interdependencies. The following topics provide an overview of the types of configuration files provided for standardization.

- [“Sun Match Engine Standardization File Types” on page 10](#)
- [“Sun Match Engine Internationalization” on page 11](#)

## Sun Match Engine Standardization File Types

Several different types of configuration files are included with the Sun Match Engine, each providing specific information to help the engine standardize and match data according to requirements. Several of these files are common to all supported nationalities, but a small subset is specific to each.

- **Category Files** - The Sun Match Engine uses category files when processing person or business names. These files list common values for certain types of data, such as titles, suffixes, and nicknames for person names or industries and organizations for business names. Category files also define standardized versions of each term or classify the terms into different categories, and some files perform both functions. When processing address files, category files named “clues files” are used.

- **Clues Files** - The Sun Match Engine uses clues files when processing address data types. These files list general terms used in street address fields, define standardized versions of each term, and classify the terms into various component types using predefined address tokens. These files are used by the standardization engine to determine how to parse a street address into its various components. Clues files provide clues in the form of tokens to help the engine recognize the component type of certain values in the input fields.
- **Constants Files** - The Sun Match Engine refers to constants files for information about the standardization files, such as the maximum length of the files. For the address data type, the constants file also describes input and output field lengths.
- **Patterns Files** - The patterns files specify how incoming data should be interpreted for standardization based on the format, or pattern, of the data. These files are used only for processing data contained in free-form text fields that must be parsed prior to matching (such as street address fields or business names). Patterns files list possible input data patterns, which are encoded in the form of tokens. Each token signifies a specific component of the free-form text field. For example, in a street address field, the house number is identified by one token, the street name by another, and so on. Patterns files also define the format of the output fields for each input pattern.
- **Key Type Files** - For business name processing, the Sun Match Engine refers to a number of key type files for processing information. These files generally define standard versions of terms commonly found in business names and some classify these terms into various components or industries. These files are used by the standardization engine to determine how to parse a business name into its different components and to recognize the component type of certain values in the input fields.
- **Reference Files** - Reference files define general terms that appear in input fields for each data type. Some reference files define terms to ignore and some define terms that indicate the business name is continuing. For example, in business name processing “and” is defined as a joining term. This helps the standardization engine to recognize that the primary business name in “Martin and Sons, Inc.” is “Martin and Sons” instead of just “Martin”. Reference files can also define characters to be ignored by the standardization engine.

## Sun Match Engine Internationalization

By default, the Sun Match Engine supports addresses and names originating from Australia, France, Great Britain, and the United States. Each national domain uses a set of common standardization files and a smaller set of unique, domain-specific files to account for international differences in address formats, names, and so on. You can process with your data using the standardization files for a single domain or you can use multiple domains depending on how the Match Field file is configured.

# Sun Match Engine Matching Configuration

The matching configuration files define how the Sun Match Engine processes records to assign matching probability weights, allowing the master index application to identify matches, potential duplicates, and non-matches. These files consist of two configurable files, the match configuration file and the match constants file. Together these files define additional logic for the Sun Match Engine to use when determining the matching probability between two records. A third file, the internal match constants file, is read-only and used internally by the match engine. It defines each comparison function and the comparison options.

The matching configuration files are very flexible, allowing you to customize the matching logic according to the type of data stored in the master index application and for the record matching requirements of your business. The matching configuration files are stored in the master index project and appear as nodes in the Match Engine node of the project. The Sun Match Engine typically standardizes the data prior to matching, so the match process is performed against the standardized data.

The matching configuration files for the Sun Match Engine must follow certain rules for formatting and interdependencies. The following topics provide an overview of the two matching configuration files provided, the architecture of those files, and formatting descriptions. They also include an overview of comparison functions used in the match configuration file.

- [“The Sun Match Engine Match Configuration File” on page 12](#)
- [“The Match Constants File” on page 18](#)

## The Sun Match Engine Match Configuration File

The match configuration file, `matchConfigFile.cfg`, contains the matching logic for each field on which matching is performed. By default, this file defines the matching logic for the three primary data types (person names, business names, and addresses), and can also handle generic data types, such as dates, numbers, social security numbers, and characters.

The match configuration file defines matching logic for each field on which matching is performed. The Sun Match Engine provides several comparison functions that you can call in this file to fine-tune the match process. Comparison functions contain the logic to compare different types of data in very specific ways in order to arrive at a match weight for each field. These functions allow you to define how matching is performed for different data types and can be used in conjunction with either matching and unmatching probabilities or agreement and disagreement weight ranges for each field. This file also defines how to handle missing fields.

The following topics describe the format of the configuration file and provide an overview of the predefined comparison functions:

- [“Sun Match Engine Match Configuration File Format” on page 13](#)
- [“Sun Match Engine Matching Comparison Functions” on page 15](#)

## Sun Match Engine Match Configuration File Format

The match configuration file is divided into two sections. The first section consists of one line that indicates the matching probability type. The second section consists of the matching rules to use for each match field.

### Match Configuration File Sample

Following is an excerpt from the default match configuration file. This excerpt illustrates the components that are described in the following sections.

```

ProbabilityType          1

FirstName                15 0 uf  0.99 0.001 15 -5
LastName                 15 0 ul  0.99 0.001 15 -5
String                   25 0 ua  0.99 0.001 10 -5
DateDays                 20 0 dD  0.99 0.001 10 -10 y 15 30
DateMonths               20 0 dM  0.99 0.001 10 -10 n
DateHours                 20 0 dH  0.99 0.001 10 -10 y 30 60
DateMinutes              20 0 dm  0.99 0.001 10 -10 y 300 600
DateSeconds              20 0 ds  0.99 0.001 10 -10 y 75 60
Numeric                  15 0 n   0.99 0.001 10 -10 y 8
Integer                  15 0 nI  0.99 0.001 10 -10 n
Real                     15 0 nR  0.99 0.001 10 -10 n
Char                      1 0 c   0.99 0.001 5 -5
pro                       15 0 p   0.99 0.001 10 -10 20 5 5

```

### Probability Type

The first line of the match configuration file defines the probability type to use for matching. Specify “0” (zero) to use m-probabilities and u-probabilities to determine a field’s match weight; specify “1” (one) to use agreement and disagreement weight ranges. If the probability type is set to use agreement and disagreement weight ranges, the **m-prob** and **u-prob** columns in the matching rules section are ignored. Likewise, if the probability type is set to use m-probabilities and u-probabilities, the **agreement-weight** and **disagreement-weight** columns in the matching rules section are ignored. The default is to use agreement and disagreement weight ranges because they are more intuitive.

### Matching Rules

The section after the first line of the match configuration file contains match field rows, with each row defining how a certain data type or field will be matched. The syntax for this section is:

match-type size null-field function m-prob u-prob agreement-weight  
disagreement-weight parameters

Table 1 describes each element in a match field row.

TABLE 1 Match Configuration File Columns

Column Number	Column Name	Description
1	match-type	A value that indicates to the Sun Match Engine how each field should be weighted. Each field included in the match string (the MatchingConfig section of the Match Field file) must have a match type corresponding to a value in this column.
2	size	The number of characters in the field on which matching is performed, beginning with the first character. For example, to match on only the first four characters in a 10-digit field, the value of this column should be "4".
3	null-field	<p>An index that specifies how to calculate the total weight for null fields or fields that only contain spaces. You can specify any of the following values:</p> <ul style="list-style-type: none"> <li>■ <b>0</b> - (zero) If one or both fields are empty, the weight used for the field is 0 (zero).</li> <li>■ <b>1</b> - (one) If both fields are empty, the agreement weight is used; if only one field is empty, the disagreement weight is used.</li> <li>■ <b>a#</b> - An "a" followed by a number specifies to use the agreement weight if both fields are empty. The agreement weight is divided by the number following the "a" to obtain the match weight for that field. If no number is specified, the default is "2". You can specify any number from 1 through 10.</li> <li>■ <b>d#</b> - A "d" followed by a number specifies to use the disagreement weight if only one field is empty. The disagreement weight is divided by the number following the "d" to obtain the match weight for the field. If no number is specified, the default is "2". You can specify any number from 1 through 10.</li> </ul> <p><b>Note</b> - In the above descriptions, the agreement and disagreement weights are either specified in this file or calculated using a logarithmic formula based on the m and u-probabilities (depending on the probability type).</p>

TABLE 1 Match Configuration File Columns (Continued)

Column Number	Column Name	Description
4	function	The type of comparison to perform when weighting the field. For information about the available comparison functions, see <a href="#">“Match Configuration Comparison Functions for Sun Match Engine (Repository)”</a> on page 93.
5	m-prob	The initial probability that the specified field in two records will match if the records match. The probability is a double value between 0 and 1, and can have up to 16 decimal points.
6	u-prob	The initial probability that the specified field in two records will match if the records do not match. The probability is a double value between 0 and 1, and can have up to 16 decimal points.
7	agreement-weight	The matching weight to be assigned to a field given that the fields match between two records. This number can be between 0 and 100 and can have up to 16 decimal points. It represents the maximum match weight for a field.
8	disagreement-weight	The matching weight to be assigned to a field given that the fields do not match between two records. This number can be between 0 and -100 and can have up to 16 decimal points. It represents the minimum match weight for a field.
9	parameters	The parameters that correspond to the comparison function specified in column 4. Some comparison functions do not take any parameters and some take multiple parameters. For additional information about parameters, see <a href="#">“Match Configuration Comparison Functions for Sun Match Engine (Repository)”</a> on page 93.

## Sun Match Engine Matching Comparison Functions

Match field comparison functions, or *comparators*, compare the values of a field in two records to determine whether the fields match. The fields are then assigned a matching weight based on the results of the comparison function. You can use several different types of comparison functions in the match configuration file to define how the Sun Match Engine should match the fields in the match string. The Sun Match Engine provides several options to use with each function.

The following table summarizes each comparison function. A complete reference of the comparison functions and their parameters is included in [“Match Configuration Comparison Functions for Sun Match Engine \(Repository\)”](#) on page 93.

TABLE 2 Comparison Functions

Comparison Function	Name	Description
b1	Bigram String Comparator	Based on the <i>Bigram</i> algorithm, this function compares two strings using all combinations of two consecutive characters and returns the total number of combinations that are the same.
b2	Advanced Bigram String Comparator	Similar to the standard Bigram comparison function (b1), but allows for character transpositions.
u	Generic String Comparator	Based on the <i>Jaro</i> algorithm, this function compares two strings taking into account uncertainty factors, such as string length, transpositions, and characters in common.
ua	Advanced Generic String Comparator	Based on the <i>Jaro</i> algorithm with variants of Winkler/Lynch and McLaughlin, this function is similar to the generic string comparator (u), but increases the agreement weight if the initial characters of each string are exact matches. This comparison function takes into account key punch and visual memory errors.
uf	Simplified String Comparator - FirstName	Based on the generic string comparator (u), this function is designed to specifically weight first name values. The string is analyzed and the weight adjusted based on statistical data.
ul	Simplified String Comparator - LastName	Based on the generic string comparator (u), this function is designed to specifically weight last name values. The string is analyzed and the weight adjusted based on statistical data.
un	Simplified String Comparator - House Numbers	Based on the generic string comparator (u), this function is designed to specifically weight house number values. The string is analyzed and the weight adjusted based on statistical data.
us	Simplified String Comparator	A custom string comparator that compares two strings taking into account such uncertainty factors as string length, transpositions, key punch errors, and visual memory errors. Unlike the generic string comparator (“u”), this function handles diacritical marks. This function also improves processing speed.
usu	Language-specific String Comparator	A custom string comparator similar to the “us” comparator with the exception that it is based in Unicode to support multiple languages and alphabets. This comparator takes one parameter indicating the language to use.
c	Exact char-by-char Comparator	Compares string fields character by character. Each character must match in order for an agreement weight to be assigned.



TABLE 2 Comparison Functions (Continued)

Comparison Function	Name	Description
n	Generic Number Comparator	Compares numeric fields using a relative distance value to determine the match weight. As the difference between the two fields increases, the match weight decreases. Once the difference is beyond the relative distance, a disagreement weight is assigned. This comparator takes two parameters; the first indicates whether to use a relative distance or direct string comparison, and the second indicates the relative distance to use.
nI	Integer Comparator	Compares integer fields using a relative distance comparison. This comparison function is based on the generic number comparator (n), and accepts the same parameters.
nR	Real Number Comparator	Compares fields containing real numbers using a relative distance comparison. This comparison function is based on the generic number comparator (n), and accepts the same parameters.
nS	Alphanumeric Comparator	Compares social security numbers or other unique identifiers, taking into account any of these parameters: <ul style="list-style-type: none"> <li>▪ Field length</li> <li>▪ Character types</li> <li>▪ Invalid values</li> </ul>
dY	Date Comparator - Year only	Compares year values using relative distance values prior to and following the given year to determine the match weight. As the difference between the two fields increases, the match weight decreases. Once the difference is beyond the relative distance, a disagreement weight is assigned. The date comparison functions handle Gregorian years. This comparator takes up to three parameters; the first indicates whether to use a relative distance or direct string comparison, and the second and third indicate the relative distance before and after.
dM	Date Comparator - Month-Year	Compares the month and year using a relative distance as described above for the year comparison function (dY).
dD	Date Comparator - Day-Month-Year	Compares the day, month, and year using a relative distance as described above for the year comparison function (dY).
dH	Date Comparator - Hour-Day-Month-Year	Compares the hour, day, month, and year using a relative distance as described above for the year comparison function (dY).
dm	Date Comparator - Min-Hour-Day-Month-Year	Compares the minute, hour, day, month, and year using a relative distance as described above for the year comparison function (dY).
ds	Date Comparator - Sec-Min-Hour-Day-Month-Year	Compares the second, minute, hour, day, month, and year using a relative distance as described above for the year comparison function (dY).

TABLE 2 Comparison Functions (Continued)

Comparison Function	Name	Description
p	Prorated Comparator	Prorates the disagreement weight for a date or numeric field based on values you specify. Differences greater than the amount you specify receive the full disagreement weight. This comparator takes three parameters indicating the relative distance and the agreement and disagreement ranges.

## The Match Constants File

The match constants file, `matchConstants.cfg`, defines certain configurable constants used by the match engine. This file includes four parameters, but currently only the first parameter, `nFields`, is used. This parameter defines the maximum number of fields being used for matching. This must be equal to or greater than the number of fields defined in the `match-columns` element of the Match Field file. The match constants file defines the following constants for the match engine.

- **nFields** - This constant defines the maximum number of different matching fields. You can enter any integer, but this number must be equal to or greater than the number of fields defined in the match-columns elements in the Match Field file.
- **maxFreqTableSize** - This constant is only used when frequency tables are used. This is not currently available and this constant is ignored.
- **maxNumberTables** - This constant is only used when frequency tables are used. This is not currently available and this constant is ignored.
- **mcls** - This constant is only used when the generic-type frequency tables are used. This is not currently available and this constant is ignored.

## Sun Match Engine and the Sun Match Engine

Implementing the Sun Match Engine with a master index application requires some customization to the Match Field file in the master index project. You can also customize the Sun Match Engine configuration files to better suit your data standardization and matching requirements.

The following topics provide information about the required customization and how the Match Field file corresponds to the configuration files.

- [“Master Index Components and the Sun Match Engine” on page 19](#)
- [“Configuring the Master Index Matching Service \(Repository\)” on page 27](#)

## Master Index Components and the Sun Match Engine

Sun Match Engine applications use the Sun Match Engine specifically for standardization and probabilistic weighting, while the master index application determines survivorship. This process relies on the logic specified in the configuration files of the master index project and of the Sun Match Engine.

The following topics provide information about how the Sun Match Engine works with master index applications to standardize data and formulate matching weights.

- “Searching and Matching in Sun Match Engine Applications (Repository)” on page 19
- “Standardization and Matching Process in Master Index Applications (Repository)” on page 19
- “The Master Index Match String (Repository)” on page 20
- “Sun Match Engine Field Identifiers” on page 20
- “Sun Match Engine Match and Standardization Types” on page 24
- “Sun Match Engine Configuration File Modifications” on page 26

### Searching and Matching in Sun Match Engine Applications (Repository)

When a new record is passed to the master index database, the master index application selects a subset of possible matches from the database. The master index application then uses the Sun Match Engine matching algorithm to assign a matching probability weight for each record in this subset (known as the candidate selection pool). To create the candidate selection pool, the master index application makes a series of query passes of the existing data, searching for matches on specific combinations of data. These combinations are defined by the blocking query, which is defined in the Candidate Select file and specified in the Threshold file.

Matching is performed on the fields included in the match string defined in the Match Field file. Each field is assigned a matching weight. The weights for each field are summed to determine the matching probability weight for the entire record (known as the *composite weight*). Before matching on some fields, such as the first name, the index might standardize the field based on information in the standardization files. You can customize how each field is weighted by modifying the match configuration file.

### Standardization and Matching Process in Master Index Applications (Repository)

The standardization and matching processes use logic that is defined by a combination of Sun Match Engine configuration files and master index configuration files. During the standardization and match processes, the following occurs.

1. The Sun Match Engine receives an incoming record.

2. The Sun Match Engine standardizes the fields specified for parsing, normalization, and phonetic encoding. These fields are defined in the StandardizationConfig section of the Match Field file and the rules for standardization are defined in the Sun Match Engine standardization configuration files.
3. The master index application queries the database for a candidate selection pool (records that are possible matches) using the blocking query specified in the Threshold file. If the blocking query uses standardized or phonetic fields, the criteria values are obtained from the database.
4. For each possible match, the master index application creates a match string (based on the match columns in the Match Field file) and sends the string to the Sun Match Engine.
5. The Sun Match Engine checks the incoming record against each possible match, producing a matching weight for each. Matching is performed using the weighting rules defined in the match configuration file.

## The Master Index Match String (Repository)

The data string that is passed to the Sun Match Engine for match processing is called the *match string* and is defined in the MatchingConfig section of the Match Field file. The Sun Match Engine configuration files, the blocking query, and the matching configuration are closely linked in the search and matching processes. The blocking query defines the select statements for creating the candidate selection pool during the matching process. The matching configuration defines the match string that is passed to the Sun Match Engine from the records in the candidate selection pool. Finally, the Sun Match Engine configuration files define how the match string is processed.

The Sun Match Engine configuration files are dependent upon the match string, and it is very important when you modify the match string to ensure that the match type you specify corresponds to the correct row in the match configuration file (`matchConfigFile.cfg`). For example, if you are using person matching and add “MaritalStatus” as a match field, you need to specify a match type for the MaritalStatus field that is listed in the first column of the match configuration file. You must also make sure that the matching logic defined in the corresponding row of the match configuration file is defined appropriately for matching on the MaritalStatus field.

## Sun Match Engine Field Identifiers

The Sun Match Engine breaks down fields into various components. For example, it breaks addresses into floor number, street number, street name, street direction, and so on. Some of these components are similar and are typically stored in the same field in the database. In the default configuration, for example, when the standardization engine finds a house number, rural route number, or PO box number, the value is stored in the HouseNumber database field. You can customize this as needed, as long as any field you specify to store a component is also included in the object structure defined for the master index application.

The Sun Match Engine uses field identifiers to determine how to process fields that are defined for normalization or parsing. The IDs are defined internally in the match engine and are referenced in the Match Field file. The field IDs you specify for each field in the Match Field file determine how that field is processed by the standardization engine. The field IDs for person names determine how each name is normalized. The field IDs for business names specify which business type key file to use for standardization. The field IDs for addresses determine which database fields store each field component and how each component is standardized.

Table 3 lists each field component generated by the Sun Match Engine along with their corresponding field IDs. You can only specify the predefined field IDs that are listed in this table.

TABLE 3 Standardization Field Identifiers

Field ID	Description
<b>Person Name Standardization Field Identifiers</b>	
FirstName	Specifies a first name field for normalization.
LastName	Specifies a last name field for normalization.
<b>Address Standardization Field Identifiers</b>	
HouseNumber	Specifies the parsed house number from a standardized address field. By default, this is stored in the <i>field_name_HouseNo</i> field (or the HouseNumber field for Sun Master Patient Index).
RuralRouteIdentif	Specifies the parsed rural route identifier from a standardized address field. By default, this is stored in the <i>field_name_HouseNo</i> field (or the HouseNumber field for Sun Master Patient Index).
BoxIdentif	Specifies the parsed PO box number from a standardized address field. By default, this is stored in the <i>field_name_HouseNo</i> field (or the HouseNumber field for Sun Master Patient Index).
MatchStreetName	Specifies the parsed and standardized street name from a standardized address field and is used internally by the match engine. If you want to store the standardized street name in the database (recommended), map this field to the street name field in the database. By default, this is stored in the <i>field_name_StName</i> field (or the StreetName field for Sun Master Patient Index).
OrigStreetName	Specifies the parsed street name from an address field. If you want to store the original street name in the database, map this field to the street name field in the database. This address component is not included in the default standardization structure, but you can add it if needed.

TABLE 3 Standardization Field Identifiers (Continued)

Field ID	Description
RuralRouteDescript	Specifies the parsed rural route description from a standardized address field. By default, this is stored in the <i>field_name_StName</i> field (or the StreetName field for Sun Master Patient Index).
BoxDescript	Specifies the PO box type from a standardized address field. By default, this is stored in the <i>field_name_StName</i> field (or the StreetName field for Sun Master Patient Index).
PropDesPrefDirection	Specifies the parsed property direction from a standardized address field. This field ID handles cases where the direction is a prefix to the property description. By default, this is stored in the <i>field_name_StDir</i> field (or the StreetDir field for Sun Master Patient Index).
PropDesSufDirection	Specifies the parsed property direction from a standardized address field. This field ID handles cases where the direction is a suffix to the property description. By default, this is stored in the <i>field_name_StDir</i> field (or the StreetDir field for Sun Master Patient Index).
StreetNamePrefDirection	Specifies the parsed street direction from a standardized address field. This field ID handles cases where the direction is a prefix to the street name. By default, this is stored in the <i>field_name_StDir</i> field (or the StreetDir field for Sun Master Patient Index).
StreetNameSufDirection	Specifies the parsed street direction from a standardized address field. This field ID handles cases where the direction is a suffix to the street name. By default, this is stored in the <i>field_name_StDir</i> field (or the StreetDir field for Sun Master Patient Index).
StreetNameSufType	Specifies the parsed street type from a standardized address field. This field ID handles cases where the street type is a suffix to the street name. By default, this is stored in the <i>field_name_StType</i> field (or the StreetType field for Sun Master Patient Index).
StreetNamePrefType	Specifies the parsed street type from a standardized address field. This field ID handles cases where the street type is a prefix to the street name. By default, this is stored in the <i>field_name_StType</i> field (or the StreetType field for Sun Master Patient Index).
PropDesSufType	Specifies the parsed property type from a standardized address field. This field ID handles cases where the street type is a suffix to the property description. By default, this is stored in the <i>field_name_StType</i> field (or the StreetType field for Sun Master Patient Index).
PropDesPrefType	Specifies the parsed property type from a standardized address field. This field ID handles cases where the street type is a prefix to the property description. By default, this is stored in the <i>field_name_StType</i> field (or the StreetType field for Sun Master Patient Index).

TABLE 3 Standardization Field Identifiers (Continued)

Field ID	Description
HouseNumPrefix	Specifies the parsed house number prefix from a standardized address field (such as the “A” in “A 1587 4th Street”). This address component is not included in the default standardization structure, but you can add it if needed.
SecondHouseNumberPrefix	Specifies the parsed <i>second</i> house number prefix from a standardized address field (such as “25” in “25 319 10th Ave.”). This address component is not included in the default standardization structure, but you can add it if needed.
SecondHouseNumber	Specifies the parsed <i>second</i> house number prefix from a standardized address field. This address component is not included in the default standardization structure, but you can add it if needed.
HouseNumSuffix	Specifies the parsed house number suffix from a standardized address field. This address component is not included in the default standardization structure, but you can add it if needed.
OrigSecondStreetName	Specifies the parsed <i>second</i> street name from a standardized address field (for example, an address might include a cross-street or a thoroughfare and dependent thoroughfare). This address component is not included in the default standardization structure, but you can add it if needed.
SecondStreetNameSufDirection	Specifies the parsed <i>second</i> street direction from a standardized address field. This address component is not included in the default standardization structure, but you can add it if needed.
SecondStreetNameSufType	Specifies the parsed <i>second</i> street type from a standardized address field. This address component is not included in the default standardization structure, but you can add it if needed.
StreetNameExtensionIndex	Specifies the parsed street name extension from a standardized address field. This address component is not included in the default standardization structure, but you can add it if needed.
WithinStructDescript	Specifies the parsed internal descriptor (such as “Floor”) from a standardized address field. This address component is not included in the default standardization structure, but you can add it if needed.
WithinStructIdentif	Specifies the parsed internal identifier (such as a floor number) from a standardized address field. This address component is not included in the default standardization structure, but you can add it if needed.
OrigPropertyName	Specifies the parsed original property name (such as the name of a complex or business park) from a standardized address field. This address component is not included in the default standardization structure, but you can add it if needed.

TABLE 3 Standardization Field Identifiers (Continued)

Field ID	Description
MatchPropertyName	Specifies the parsed match property name from a standardized address field and is used internally by the match engine for blocking and phonetic encoding. This address component is not included in the default standardization structure, but you can add it if needed.
CenterDescript	Specifies the parsed structure description from a standardized address field. This address component is not included in the default standardization structure, but you can add it if needed.
CenterIdentif	Specifies the parsed structure identifier from a standardized address field. This address component is not included in the default standardization structure, but you can add it if needed.
ExtraInfo	Specifies any extra information that was not included in any of the other parsed components. This address component is not included in the default standardization structure, but you can add it if needed.
<b>Business Name Standardization Field Identifiers</b>	
PrimaryName	Specifies the field containing the parsed name in a free-form text business name field.
OrgTypeKeyword	Specifies the field containing the parsed organization type in a free-form text business name field.
AssocTypeKeyword	Specifies the field containing the parsed association type in a free-form text business name field.
IndustrySectorList	Specifies the field containing the parsed industry sector in a free-form text business name field.
IndustryTypeKeyword	Specifies the field containing the parsed industry type in a free-form text business name field (industry type is a subset of the sector).
AliasList	Specifies the field containing the parsed alias in a free-form text business name field.
Url	Specifies the field containing the parsed URL in a free-form text business name field.

## Sun Match Engine Match and Standardization Types

Indicators are used in the Match Field file to reference the type of matching and standardization to perform on each field. You must specify one of these indicators, called *match types* and *standardization types*, for the fields you define for standardization or matching. The match types correspond to the match types listed in the first column of the match configuration file (`matchConfigFile.cfg`). The standardization types are defined internally in the match engine. The Sun Match Engine uses these types to determine how to process each field.



Table 4 lists the default standardization types; Table 5 lists the default match types. You can modify the match type names but not the standardization type names. For more information about match and standardization types, see “[Master Index Match Types and Field Names \(Repository\)](#)” in *Understanding Sun Master Index Processing (Repository)*. Note that the match types you can specify in the Match Field file (listed in Table 5) are not the same values you specify for the Match Type field drop-down list in the wizard when you create the master index application.

TABLE 4 Standardization Types

This indicator ...	processes this data type ...
Address	Free-form street address fields.
PersonName	Pre-parsed name fields (including any first, middle, last, or alias names).
BusinessName	Free-form business names.

The standardization types listed above correspond to the three categories of match types listed below. You can also specify miscellaneous match types, which do not correspond to any standardization types.

TABLE 5 Match Types

This indicator ...	processes this data type ...
<b>Business Name Match Types</b>	
PrimaryName	The parsed name field of a business name.
OrgTypeKeyword	The parsed organization type field of a business name.
AssocTypeKeyword	The parsed association type field of a business name.
AliasList	The parsed alias type field of a business name.
IndustrySectorList	The parsed industry sector field of a business name.
IndustryTypeKeyword	The parsed industry type field of a business name.
Url	The parsed URL field of a business name.
<b>Address Match Types</b>	
StreetName	The parsed street name field of a street address.
HouseNumber	The parsed house number field of a street address.
StreetDir	The parsed street direction field of a street address.
StreetType	The parsed street type field of a street address.

TABLE 5 Match Types (Continued)

This indicator ...	processes this data type ...
<b>Person Name Match Types</b>	
FirstName	A first name field, including middle name, alias first name, and alias middle name fields.
LastName	A last name field, including alias last name fields.
<b>Date Match Types</b>	
DateDays	The day, month, and year of a date field.
DateMonths	The month and year of a date field.
DateHours	The hour, day, month, and year of a date field.
DateMinutes	The minute, hour, day, month, and year of a date field.
DateSeconds	The seconds, minute, hour, day, month, and year of a date field.
<b>Miscellaneous Match Types</b>	
String	A generic string field.
Numeric	A numeric field.
Integer	A field containing integers.
Real	A field containing real numbers.
SSN	A field containing a social security number.
Char	A field containing a single character.
pro	Any field on which you want the Sun Match Engine to use prorated weights.
Exac	Any field you want the Sun Match Engine to match character for character.

## Sun Match Engine Configuration File Modifications

The Sun Match Engine configuration files are designed to perform very specific functions in the standardization and match processes. These files should only be modified by personnel with an understanding of the Sun Match Engine and an understanding of the data integrity requirements of your organization. Modifications to both the master index configuration files and the Sun Match Engine configuration files should be made while the master index application is in the preproduction stages. Modifying the files after master index application has moved into production might cause variances in matching weights and data processing.

The most common modifications to the Sun Match Engine configuration files are generally in the match configuration file, where you can fine-tune the weighting process. This file defines probabilities used by the algorithm to determine a matching probability weight for each match field. You can use the match comparison functions provided by the Sun Match Engine to

fine-tune the matching logic in this file. Another common modification is inserting additional names or terms into category files, such as the first name category file (`personFirstName*.dat`).

Depending on your data requirements, you might need to modify additional standardization files. Some of the patterns files (most notably the address patterns files) are very complex and should only be modified by personnel who thoroughly understand the defined patterns and tokens. If you modify standardization files, make sure you modify them for each national domain specified in the Match Field file.

## Configuring the Master Index Matching Service (Repository)

To configure a Sun Match Engine application for specific data types and for the Sun Match Engine, you must customize the Matching Service by modifying the Match Field file in the master index project. Configuring the matching service consists of the following four tasks.

- [“Master Index Standardization Configuration \(Repository\)” on page 27](#)
- [“Master Index Match String Configuration \(Repository\)” on page 30](#)
- [“Match and Standardization Engine Configuration” on page 30](#)
- [“Master Index Phonetic Encoder Configuration \(Repository\)” on page 31](#)

### Master Index Standardization Configuration (Repository)

The *StandardizationConfig* section of the Match Field file determines which fields are normalized, parsed, or phonetically encoded and defines the nationality of the data being processed. The standardization section includes the following structures.

- [“Normalization Structures” on page 27](#)
- [“Standardization Structures \(Parsing and Normalization\)” on page 28](#)
- [“Phonetic Encoding Structures” on page 29](#)

The *StandardizationConfig* section defines fields that will be normalized, fields that will be parsed and normalized, and fields that will be phonetically encoded. The standardization types you specify in this section correspond to the match configuration file; the field IDs you can specify are listed in [Table 3](#).

### Normalization Structures

The normalization structure defines fields that are already parsed, but need to be normalized. It also tells the Sun Match Engine where to place the normalized data in the object structure. Matching on any of these fields is determined by the match string and the logic is defined in the match configuration file.

Of the three data types processed by the Sun Match Engine, only the person name data type is expected to provide information in fields that are already parsed; that is, the first, last, and middle names appear in separate fields, as do the suffix, title, and so on. The person standardization files define logic for normalizing person name fields. By default, only the names you specify for matching in the wizard are defined for normalization. You can define normalization for additional name fields, such as maiden name, spouse's name, and so on. For each normalization structure, you must specify the national domains for the data you are processing.

### **Defining New Fields for Normalization**

The fields you define for normalization in the Match Field file can include any name fields. If you define normalization for fields that are not currently defined for normalization in the Match Field file, make the following additional changes.

1. In the Match Field file, define the normalization structure, using the appropriate standardization type (PersonName), domain selector, and field IDs (FirstName, MiddleName, or LastName).
2. Add the new fields that will store the normalized field value to the appropriate objects in the Object Definition file.
3. If any of the normalized fields are to be used for blocking, modify the Candidate Select file by adding the new fields to the blocking query.
4. Regenerate the master index application in NetBeans to include the new fields in the database creation script, the outbound Object Type Definition (OTD), and the method OTD.
5. To specify that the new normalized fields be used for matching, do the following:
  - a. Determine the match type or the match comparison function you want to use to match the normalized data, and modify the match configuration file (`matchConfigFile.cfg`) if needed.
  - b. Add the new normalized field to the *match-columns* element of the MatchingConfig section of the Match Field file, making sure to use the appropriate match type from the match configuration file.

### **Standardization Structures (Parsing and Normalization)**

The fields that must be parsed, and possibly normalized, are defined in a standardization structure in the StandardizationConfig section of the Match Field file. The standardization structure tells the Sun Match Engine where to place the standardized information extracted from the parsed fields. The target fields you specify for standardization facilitate searching by the parsed values. Matching on any of these fields is determined by the match string and the logic is defined in the match configuration file.

The Sun Match Engine expects business names and street address information in free-form text fields that must be parsed and normalized prior to matching. The logic for parsing and

normalizing street address information is contained in the address standardization files; the logic for parsing and normalizing business names is contained in the business standardization files. You can customize the standardization of these data types by modifying the appropriate patterns file. For each standardization structure, you must specify the national domains for the data being processed.

### Defining New Fields for Standardization

The fields you define for standardization in the Match Field file can include any street address or business name field. Perform the following steps if you need to define one of these field types for standardization.

1. If necessary, modify the patterns file for the type of data you are standardizing.  
You can define new input and output patterns or modify existing ones.
2. Define the standardization structure, using the appropriate standardization type (BusinessName or Address), domain selector, and field IDs (described in [Table 3](#)).
3. Add the new fields that will store the parsed or normalized data to the appropriate objects in the Object Definition file.
4. If any of the parsed or normalized fields are to be used for blocking, modify the Candidate Select file by adding the new fields to the blocking query.
5. Regenerate the master index application in NetBeans to include the new fields in the database creation script, the outbound Object Type Definition (OTD), and the method OTD.
6. To specify that the new standardized fields be used for matching, do the following:
  - a. Determine the match type or the match comparison function you want to use to match the parsed data, and modify the match configuration file (`matchConfigFile.cfg`) if needed.
  - b. Add the new standardized field to the *match-columns* element of the MatchingConfig section of the Match Field file, making sure to use the appropriate match type from the match configuration file.

### Phonetic Encoding Structures

The fields to be phonetically encoded are defined in a phonetic encoding structure in the StandardizationConfig section of the Match Field file. The phonetic encoding structure tells the Sun Match Engine where to place the phonetic data created from the fields that are encoded. You can define any field in the object structure for phonetic encoding.

### Defining New Fields for Phonetic Encoding

The fields you define for phonetic encoding in the Match Field file can include any field.

1. Determine the type of phonetic encoder to use to convert the field.  
You can use any of the encoders described in [Table 7](#).
2. Define the phonetic encoding structure, using the appropriate encoders.
3. Add the new fields that will store the phonetic values to the appropriate objects in the Object Definition file.
4. If any the phonetic fields are to be used for blocking, modify the Candidate Select file by adding the new fields to the blocking query.
5. Regenerate the master index application in NetBeans to include the new fields in the database creation script, the outbound OTD, and the method OTD.

## Master Index Match String Configuration (Repository)

The *MatchingConfig* section of the Match Field file determines which fields are passed to the Sun Match Engine for matching (the match string). If you are matching on fields parsed from a free-form text field, define each individual parsed field you want to use for matching. The default fields listed in the *MatchingConfig* section depend on the fields you specified for matching in the wizard (for Sun Master Patient Index, the default fields are *FirstName*, *LastName*, *DOB*, *Gender*, and *SSN*).

The match types you can use for each field in this section are defined in the first column of the match configuration file. Make sure the match type you specify has the correct matching logic defined in the match configuration file.

## Match and Standardization Engine Configuration

The *MEFAConfig* section of the Match Field file defines which standardization and match engines will be used by the master index application. By default, the master index application is already configured to use the Sun Match Engine for matching and standardization. For more information, see [Understanding Sun Master Index Configuration Options \(Repository\)](#).

[Table 6](#) lists the elements in the Match Field file that define the match and standardization engine, along with the appropriate values for the Sun Match Engine.

TABLE 6 Sun Match Engine Standardization and Match Classes

Match Field File Element	Sun Match Engine Value
standardizer-api	com.st.eindex.matching.adapter.SbmeStandardizerAdapter
standardizer-config	com.st.eindex.matching.adapter.SbmeStandardizerAdapter Config
matcher-api	com.st.eindex.matching.adapter.SbmeMatcherAdapter

TABLE 6 Sun Match Engine Standardization and Match Classes (Continued)

Match Field File Element	Sun Match Engine Value
matcher-config	com.stc.eindex.matching.adapter.SbmeMatcherAdapter Config

## Master Index Phonetic Encoder Configuration (Repository)

The Sun Match Engine supports several phonetic encoders, which are defined in the *PhoneticEncodersConfig* section of the Match Field file. Any encoders specified in the phonetic encoding structures (see “[Phonetic Encoding Structures](#)” on page 29) must also be defined in the *PhoneticEncodersConfig* section. The classes for the encoders are listed in [Table 7](#).

- **Soundex** - This algorithm is an industry standard for phonetically encoding first names.
- **French Soundex** - This algorithm is based on the Soundex algorithm, but is customized for French characters and names.
- **Refined Soundex** - This algorithm is similar to the Soundex algorithm, but is optimized for spell checking.
- **NYSIIS** - This algorithm is an industry standard for phonetically encoding last names.
- **Metaphone** - This algorithm is similar to the Soundex algorithm, but is better at identifying words that sound similar. This encoder is limited to encoding a single word in ASCII format containing only characters in the A - Z range. No punctuation or numbers can be in the input string.
- **Double Metaphone** - This algorithm is an improvement on the Metaphone algorithm, at times returning two encodings for a word that could have multiple pronunciations.

TABLE 7 Phonetic Encoder Classes for the Sun Match Engine

Encoder	Java Class
Soundex	com.stc.eindex.phonetic.impl.Soundex
NYSIIS	com.stc.eindex.phonetic.impl.NYSIIS
Metaphone	com.stc.eindex.phonetic.impl.Metaphone
Double Metaphone	com.stc.eindex.phonetic.impl.DoubleMetaphone
Refined Soundex	com.stc.eindex.phonetic.impl.RefinedSoundex
French Soundex	com.stc.eindex.phonetic.impl.SoundexFR

# Sun Match Engine Person Data Type Configuration

Processing person data involves normalizing and phonetically encoding certain fields prior to matching. The following topics describe the default configuration files that define person processing logic and provide instructions for modifying the Match Field file for processing person data. The information presented in this in these topics is especially pertinent to the Sun Master Patient Index application.

- [“Sun Match Engine Person Matching Overview”](#) on page 32
- [“Sun Match Engine Match Configuration for Person Data”](#) on page 34
- [“Sun Match Engine Person Data Standardization Files”](#) on page 34
- [“Configuring the Sun Match Engine Standardization Files for Person Data”](#) on page 43
- [“Configuring the Master Index Matching Service for Person Data \(Repository\)”](#) on page 43

## Sun Match Engine Person Matching Overview

Matching on the person data type includes standardizing and matching a person’s demographic information. The Sun Match Engine can create normalized and phonetic values for person data. Several configuration files designed specifically to handle person data are included to provide additional logic for the standardization and phonetic encoding process. The Sun Match Engine can phonetically encode any field you specify, with some modification to the standardization files. It can also match on any field, as long as the match type for the field is defined in the match configuration file (`matchConfigFile.cfg`).

In addition, when storing person information, you might want to standardize addresses to enable searching against address information. This requires working with the address configuration files described in [“Sun Match Engine Address Data Type Configuration”](#) on page 48.

The following topic provides information about the fields used in person data matching and the fields added to the object structure.

- [“Sun Match Engine Person Data Processing Fields”](#) on page 32

## Sun Match Engine Person Data Processing Fields

When matching on person data, not all fields in a record need to be processed by the Sun Match Engine. The match engine only needs to process fields that must be parsed, normalized, or phonetically converted, and the fields against which matching is performed. These fields are defined in the Match Field file and processing logic for each field is defined in the Sun Match Engine standardization and matching configuration files.



## Person Data Match String Fields

The match string processed by the Sun Match Engine is defined by the match fields specified in the Match Field file. The match engine can process any combination of fields you specify for matching. By default, the match configuration file (`matchConfigFile.cfg`) includes rows specifically for matching on first name, last name, social security number, and dates (such as a date of birth). It also includes a row for matching a single character, such as might be the case in a gender field. You can use any of the existing rows for matching or you can add rows for the fields you want to match. Any field for which you specify a match type in the wizard is added to the match string.

## Person Data Standardized Fields

The Sun Match Engine expects person data to be provided in separate fields within a single record, meaning that no parsing is required of the name fields prior to normalization. Typically, only first and last names are normalized and phonetically encoded when standardizing person data, but the match engine can normalize and phonetically encode any field you choose.

## Person Data Object Structure

The fields you specify for person name matching in the wizard are automatically defined for standardization and phonetic encoding. If you specify the appropriate match types in the wizard, the following fields are automatically added to the object structure and database creation script.

- `field_name_Std`
- `field_name_Phon`

where `field_name` is the name of the field for which you specified person name matching. For example, if you specify the `PersonFirstName` match type for the `FirstName` field, two fields, `FirstName_Std` and `FirstName_Phon`, are automatically added to the structure. You can also add these fields manually if you do not specify match types in the wizard. If you store additional names in the database, such as alias names, maiden names, parent names, and so on, you can modify the phonetic structure to phonetically encode those names as well.

---

**Note** – The object structure for Sun Master Patient Index uses a slightly different naming convention. For the names of the fields defined for Sun Master Patient Index, refer to the *Sun Patient Master Index Configuration Options*.

---

## Sun Match Engine Match Configuration for Person Data

The default match configuration file, `matchConfigFile.cfg`, defines several match types for the kinds of data typically included in a person master index application. You can customize the existing match types or create new match types for the data being processed. The following match types are typical for matching on person data.

- 
- |             |          |
|-------------|----------|
| ▪ FirstName | ▪ Real   |
| ▪ LastName  | ▪ SSN    |
| ▪ String    | ▪ Gender |
| ▪ Date      | ▪ pro    |
| ▪ Numeric   | ▪ Exac   |
| ▪ Integer   |          |
- 

This file appears under the Match Engine node of the master index project. For more information about the comparison functions used for each match type and how the weights are tuned, see [“Customizing the Match Configuration” on page 88](#) and [“Match Configuration Comparison Functions for Sun Match Engine \(Repository\)” on page 93](#).

## Sun Match Engine Person Data Standardization Files

Several configuration files are used to define standardization logic for the Sun Match Engine. You can customize any of the configuration files described in this section to fit your processing and standardization requirements for person data. There are two types of standardization files for person data: common and domain-specific. The common files appear under the Standardization Engine node of the master index project and are used for all national domains; the domain-specific files appear within sub-folders of the Standardization Engine node and each corresponds to a specific national domain.

The following topics provide information about each type of standardization file:

- [“Sun Match Engine Common Standardization Files for Person Data” on page 34](#)
- [“Sun Match Engine Domain-Specific Standardization Files for Person Data” on page 36](#)

### Sun Match Engine Common Standardization Files for Person Data

The standardization files described in this section are common to all national domains. These files define special characters to remove from name fields and define hyphenated first names. A patterns file is also common, but is not currently used.

## The Hyphenated Name Category File (personFirstNameDash.dat)

The hyphenated name category file defines first names that include hyphens (such as Anne-Marie) to help the Sun Match Engine recognize and process these values as first names. The file also classifies each name into a gender category. This file is used to standardize all domains except Australia, which uses the personFirstNameDashAU.dat file located in the Australia folder, and France, which uses the personFirstNameDashFR.dat file located in the France folder.

The hyphenated name category files use the following syntax:

```
name gender-class
```

You can modify or add entries in this table as needed. [Table 8](#) describes the columns in the personFirstNameDash.dat file.

**TABLE 8** Hyphenated Name Category File

Column	Description
name	A hyphenated first name.
gender-class	An indicator of the gender with which the first name corresponds. The possible values are: <ul style="list-style-type: none"> <li>■ N - The name is neutral, and can be applied to male or female first names.</li> <li>■ F - The name is used for females.</li> <li>■ M - The name is used for males.</li> </ul>

Following is an excerpt from the personFirstNameDash.dat file.

```
ANNE-MARIE      F
JEAN-NOEL       M
JEAN-MARIE      M
JEAN-BAPTISTE   M
JEAN-PIERRE     M
JEAN-YVES       M
```

## The Person Name Patterns File (personNamePatt.dat)

The person name patterns file is not currently used, but is designed to standardize free-form text name fields.

## The Special Characters Reference File (personRemoveSpecChars.dat)

The special characters reference file lists characters that might appear in person data, but that should be ignored. The Sun Match Engine removes these characters from a field before making any comparisons or before normalizing data. You can define additional characters to remove from person data by simply adding the character to the list.

An excerpt from the `personRemoveSpecChars.dat` file appears below.

```
[  
]  
{  
}  
<  
>  
/  
?  
*  
^  
#  
!
```

## Sun Match Engine Domain-Specific Standardization Files for Person Data

Most standardization files for person data are specific to each national domain. Each domain node within the Standardization node of the project includes the files defined in this section. The domain corresponding to each file is indicated at the end of the file name; for example, `personConstantsUK.cfg` and `personConstantsFR.cfg`. These domain abbreviations are indicated by an asterisk (\*) in the descriptions.

---

**Note** – You can customize these files to add entries of other nationalities or languages, including those containing diacritical marks.

---

### The Conjunction Reference File (`personConjon*.dat`)

The conjunction reference file is not currently used, but is designed to work with the person name patterns file during standardization.

### The Person Constants File (`personConstants*.cfg`)

The person constants file defines certain information about the standardization files used for processing person data, primarily the number of lines contained in each file. The number of lines specified here must be equal to or greater than the number of lines actually contained in each file. The constants file for United States data is in the Standardization node of the project and is named `personConstants.cfg`; the person constants file for the other domains is located under the domain name node.

[Table 9](#) lists and describes each parameter in the constants file. The files referenced by these parameters are described on the following pages.

TABLE 9 Person Constants File Parameters

Parameter	Description
words	The maximum number of words in a given free-form text field containing a person name. This parameter is not currently used.
conjmax	The maximum number of lines in the person conjunction reference file (personConjon*.dat).
jrsmx	The maximum number of lines in the generational suffix category file (personGenSuffix*.dat).
nickmax	The maximum number of lines in the first name category file (personFirstName*.dat).
lastmax	The maximum number of lines in the last name category file (personLastName*.dat).
premax	The maximum number of lines in the last name prefix category file (personLastNamePrefix*.dat).
titlmax	The maximum number of lines in the title category file (personTitle*.dat).
sufmax	The maximum number of lines in the occupational suffix category file (personOccupSuffix*.dat).
skpmax	The maximum number of lines in the business name reference file (businessOrRelated*.dat).
ptrnmax1	The maximum number of lines in the person patterns file (personNamePatt.dat).
twomax	The maximum number of lines in the two-character reference file for occupational suffixes (personTwo*.dat).
thremax	The maximum number of lines in the three-character reference file for occupational suffixes (personThree*.dat).
blnkmax	The maximum number of lines in the special characters reference file (personRemoveSpecChars.dat).
dashSize	The maximum number of lines in the hyphenated name category file (personFirstNameDash.dat).

### The First Name Category File (personFirstName\*.dat)

The first name category file defines standardized versions of first names and assigns a gender classification for each name. This file is used to standardize first names when comparing person names. The gender classification helps to further clarify the match. The Sun Match Engine uses this file when a first name field is defined for normalization or standardization in the Match Field file.

The syntax of this file is:

original-value standardized-form gender-class

You can modify or add entries in this table as needed. [Table 10](#) describes the columns in the `personFirstName*.dat` file.

**TABLE 10** First Name Category File

Column	Description
original-value	The original value of the first name.
standardized-form	The standardized version of the original value. A zero (0) in this field indicates that the original value is already in its standardized form.  If this column contains a name instead of a zero, that name must also be listed in a different entry as an original value with a standardized form of "0".
gender-class	An indicator of the gender with which the first name corresponds. The possible values are: <ul style="list-style-type: none"> <li>■ N – The name is neutral, and can be applied to male or female first names.</li> <li>■ F – The name is used for females.</li> <li>■ M – The name is used for males.</li> </ul>

Following is an excerpt from the `personFirstNameUS.dat` file. Certain rows contain a zero (0) for the standardized form, indicating that the name is already standard (for example, Stephen, Sterling, and Summer).

```
STEPHEN      0           M
STEPHENIE    STEPHANIE  F
STEPHIE      STEPHANIE  F
STEPHINE     STEPHANIE  F
STEPHNIE     STEPHANIE  F
STERLING     0           M
STEVE        STEPHEN    M
STEVEN       STEPHEN    M
STEVIE       STEPHEN    N
STEW         STUART     M
STEWART      STUART     M
STU          STUART     M
STUART       0           M
SU           SUSAN      F
SUE          SUSAN      F
SUHANTO      0           M
SULLIVAN     0           F
SULLY        SULLIVAN   F
SUMMER       0           F
```

## The Generational Suffix Category File (personGenSuffix\*.dat)

The generational suffix category file defines standardized versions of generational suffixes, such as Jr., III, and so on. This file is used to compare standard versions of the suffix field. You can define additional suffixes and their standardized form following the syntax below.

field-value standard-form

Table 11 describes each column of the personGenSuffix\*.dat file.

TABLE 11 Generational Suffix Category File

Column	Description
field-value	The original value of the generational suffix in the record being processed.
standard-form	The standard form of the generational suffix. A zero (0) in this column indicates that the value listed in column one is already in its standardized form.  If this column contains a suffix instead of a zero, that suffix must also be listed in a different entry as an original value with a standard form of "0".

An excerpt from the personGenSuffixUS.dat file appears below. In this excerpt, certain suffixes, such as 2ND, 3RD and JR, are already in their standardized form.

```

11          2ND
111         3RD
1V          4TH
2ND         0
3RD         0
4TH         0
FOURTH     4TH
II          2ND
III         3RD
IV          4TH
JR          0
JUNIOR     JR
SECOND     2ND
SENIOR     SR

```

## Last Name Prefix Category File (personLastNamePrefix\*.dat)

The last name prefix category file defines standardized versions of last name prefixes, such as "Van" or "Le". This file is used to standardize these prefixes prior to standardizing the last name when comparing person names. The Sun Match Engine uses this file when a last name field is defined for normalization or standardization in the Match Field file.

The syntax of this file is:

original-value standardized-form

You can modify or add entries in this table as needed. [Table 12](#) describes the columns in the `personLastNamePrefix*.dat` file.

**TABLE 12** Last Name Prefix Category File

Column	Description
original-value	The original value of the last name prefix.
standardized-form	The standardized version of the original value. A zero (0) in this field indicates that the original value is already in its standardized form.  If this column contains a prefix instead of a zero, that prefix must also be listed in a different entry as an original value with a standardized form of "0".

Following is an excerpt from the `personLastNamePrefixUS.dat` file. Some of these prefixes are already in their standardized form, such as "Los" and "Mac".

```

LOS           0
MAC           0
MC           MAC
SAINT        0
ST           SAINT
VAN          0
VAN DER      0
VANDE       VAN DER

```

## The Last Name Category File (`personLastName*.dat`)

The last name category file defines standardized versions of last names. This file is used to standardize last names when comparing person names. The Sun Match Engine uses this file when a last name field is defined for normalization or standardization in the Match Field file.

The syntax of this file is:

original-value standardized-form

You can modify or add entries in this table as needed. [Table 13](#) describes the columns in the `personLastName*.dat` file.

**TABLE 13** Last Name Category File

Column	Description
original-value	The original value of the last name.



TABLE 13 Last Name Category File (Continued)

Column	Description
standardized-form	The standardized version of the original value. A zero (0) in this field indicates that the original value is already in its standardized form.  If this column contains a name instead of a zero, that name must also be listed in a different entry as an original value with a standardized form of "0".

Following is an excerpt from the personLastNameUS.dat file.

```
FINK                0
PHINQUE            FINK
```

### The Occupational Suffix Category File (personOccupSuffix\*.dat)

The occupational suffix category file is not currently used, but is designed to work with the person name patterns file during standardization.

### The Three-Character Suffix File (personThree\*.dat)

This reference file is not currently used, but is designed to work with the person name patterns file during standardization.

### The Title Category File (personTitle\*.dat)

The title category file defines standard forms for titles and classifies each title into a gender category. For example, "Mister" is standardized to "MR" and is classified as male; "Doctor" is standardized to "DR" and is classified as gender neutral. You can add, modify, or delete entries in this file as needed. Use the following syntax.

```
original-value standardized-form gender-class
```

Table 14 describes each column of the personTitle\*.dat file.

TABLE 14 Person Title Category File

Column	Description
original-value	The original value of the title in the person name field.
standardized-form	The standardized version of the original value. A zero (0) in this field indicates that the original value is already in its standardized form.  If this column contains a title instead of a zero, that title must also be listed in a different entry as an original value with a standardized form of "0".

TABLE 14 Person Title Category File (Continued)

Column	Description
gender-class	<p>An indicator of the gender with which the title corresponds. The default values are:</p> <ul style="list-style-type: none"> <li>■ N – The title is neither male nor female.</li> <li>■ F – The title is used for females.</li> <li>■ M – The title is used for males.</li> </ul>

An excerpt from the personTitleUS.dat file appears below. In this excerpt, certain titles, such as DR, GEN, and MISS, are already in their standardized form.

CTO	0	N
DEAN	0	N
DIR	DIRECTOR	N
DIRECTOR	0	N
DOC	DR	N
DOCTOR	DR	N
DR	0	N
DRS	0	N
EMERITUS	0	N
FOUNDER	0	N
GEN	0	N
GENERAL	GEN	N
MANAGER	0	N
MGR	MANAGER	N
MISS	0	F
MISSUS	MRS	F

### The Two-Character Suffix File (personTwo\*.dat)

This reference file is not currently used, but is designed to work with the person name patterns file during standardization.

### The Business-Related Category File (businessOrRelated\*.dat)

The business-related category file is used to identify business terms in person name information. Examples of when this could occur would be when indexing both person and business names or when business information is included within a person object structure. The Sun Match Engine removes these terms for person matching. This file contains a list of common business terms that might be found in person data. You can modify this file by adding, changing, or deleting terms.

An excerpt from the businessOrRelatedUS.dat file appears below.

ACCOUNTANT  
ACCT  
ACDY  
ACRE  
ACREAGE  
ACRES  
ACS  
ACT  
AD  
ADATU  
ADM  
ADMIN  
ADMINISTRATIO  
ADMINISTRATION  
ADMINISTRATOR

## Configuring the Sun Match Engine Standardization Files for Person Data

To customize the Sun Match Engine configuration files for processing person data, you can modify any of the person data standardization files using the text editor provided in NetBeans. Before modifying the match configuration file, review the information provided in [“Sun Match Engine Standardization Configuration” on page 10](#) and [“Match Configuration Comparison Functions for Sun Match Engine \(Repository\)” on page 93](#). Make sure a thorough data analysis has been performed to determine the best fields for matching and the best comparison functions to use for each field.

Updating most standardization files is a straightforward process. Make sure to follow the syntax guidelines provided in [“Sun Match Engine Person Data Standardization Files” on page 34](#). If you add any lines to any of the standardization configuration files, be sure to adjust the corresponding parameter in the person constants file (`personConstants*.cfg`).

## Configuring the Master Index Matching Service for Person Data (Repository)

To ensure correct processing of person information, you must customize the Matching Service. This includes modifying the Match Field file to support the fields on which you want to match, to standardize the appropriate fields, and to specify the Sun Match Engine as the match and standardization engine (by default, the Sun Match Engine is already specified so this does not need to be changed). Perform the following tasks to configure the Matching Service.

- [“Configuring the Standardization Structure for Person Data \(Repository\)” on page 44](#)
- [“Configuring the Match String for Person Data \(Repository\)” on page 47](#)

When configuring the Matching Service, keep in mind the information presented in [“Configuring the Master Index Matching Service \(Repository\)” on page 27](#).

## Configuring the Standardization Structure for Person Data (Repository)

The standardization structure is configured in the StandardizationConfig section of the Match Field file, which is described in detail in [Understanding Sun Master Index Configuration Options \(Repository\)](#). To configure the required fields for normalization and phonetic encoding, modify the normalization and phonetic encoding structures in the Match Field file. The following sections provide additional guidelines and samples specific to standardizing person data.

---

**Note** – In the current configuration, the rules defined for the person data type assume the incoming data to be parsed prior to processing. Therefore, you do not need to configure fields to parse unless you want to search on address information. In that case, you must configure the address fields to parse and normalize.

---

### Person Data Normalization Structures

The fields defined for normalization for the person data type can include any name fields. By default this includes first, middle, and last name fields. Follow the instructions under “Defining Master Index Normalization Rules (Repository)” in [Configuring Sun Master Indexes \(Repository\)](#) to define fields for normalization. For the *standardization-type* element, enter **PersonName** (for more information, see [“Sun Match Engine Match and Standardization Types” on page 24](#)). For a list of field IDs to use in the *standardized-object-field-id* element, see [Table 3](#).

A sample normalization structure for person data is shown below. This sample specifies that the PersonName standardization type is used to normalize the first name, alias first name, last name, and alias last name fields. For all name fields, both United States and United Kingdom domains are defined for standardization.

```
<structures-to-normalize>
  <group standardization-type="PersonName"
    domain-selector="com.stc.eindex.matching.impl.MultiDomainSelector">
    <locale-field-name>Person.PobCountry</locale-field-name>
    <locale-maps>
      <locale-codes>
        <value>UNST</value>
        <locale>US</locale>
      </locale-codes>
      <locale-codes>
        <value>GB</value>
        <locale>UK</locale>
      </locale-codes>
    </locale-maps>
  </group>
</structures-to-normalize>
```

```

<unnormalized-source-fields>
  <source-mapping>
    <unnormalized-source-field-name>Person.FirstName
  </unnormalized-source-field-name>
  <standardized-object-field-id>FirstName
  </standardized-object-field-id>
</source-mapping>
<source-mapping>
  <unnormalized-source-field-name>Person.LastName
</unnormalized-source-field-name>
  <standardized-object-field-id>LastName
  </standardized-object-field-id>
</source-mapping>
</unnormalized-source-fields>
<normalization-targets>
  <target-mapping>
    <standardized-object-field-id>FirstName
  </standardized-object-field-id>
  <standardized-target-field-name>Person.FirstName_Std
  </standardized-target-field-name>
</target-mapping>
  <target-mapping>
    <standardized-object-field-id>LastName
  </standardized-object-field-id>
  <standardized-target-field-name>Person.LastName_Std
  </standardized-target-field-name>
</target-mapping>
</normalization-targets>
</group>
<group standardization-type="PersonName" domain-selector=
"com.stc.eindex.matching.impl.MultiDomainSelector">
  <locale-field-name>Person.PobCountry</locale-field-name>
  <locale-maps>
    <locale-codes>
      <value>UNST</value>
      <locale>US</locale>
    </locale-codes>
    <locale-codes>
      <value>GB</value>
      <locale>UK</locale>
    </locale-codes>
  </locale-maps>
  <unnormalized-source-fields>
    <source-mapping>
      <unnormalized-source-field-name>Person.Alias[*].FirstName
    </unnormalized-source-field-name>
    <standardized-object-field-id>FirstName
  </standardized-object-field-id>

```

```
</source-mapping>
<source-mapping>
  <unnormalized-source-field-name>Person.Alias[*].LastName
</unnormalized-source-field-name>
  <standardized-object-field-id>LastName
</standardized-object-field-id>
</source-mapping>
</unnormalized-source-fields>
<normalization-targets>
  <target-mapping>
    <standardized-object-field-id>FirstName
  </standardized-object-field-id>
  <standardized-target-field-name>
    Person.Alias[*].FirstName_Std
  </standardized-target-field-name>
</target-mapping>
  <target-mapping>
    <standardized-object-field-id>LastName
  </standardized-object-field-id>
  <standardized-target-field-name>
    Person.Alias[*].LastName_Std
  </standardized-target-field-name>
</target-mapping>
</normalization-targets>
</group>
</structures-to-normalize>
```

## Person Data Phonetic Encoding

When you specify a name field for person name matching in the wizard, these fields are automatically defined for phonetic encoding. You can define additional names, such as maiden names or alias names, for phonetic encoding as well. Follow the instructions under “Defining Phonetic Encoding for the Master Index (Repository)” in *Configuring Sun Master Indexes (Repository)* to define fields for phonetic encoding.

A sample of fields defined for phonetic encoding is shown below. This sample converts name and alias name fields, as well as the street name.

```
<phoneticize-fields>
  <phoneticize-field>
    <unphoneticized-source-field-name>Person.FirstName_Std
  </unphoneticized-source-field-name>
  <phoneticized-target-field-name>Person.FirstName_Phon
  </phoneticized-target-field-name>
  <encoding-type>Soundex</encoding-type>
</phoneticize-field>
  <phoneticize-field>
    <unphoneticized-source-field-name>Person.LastName_Std
```

```

    </unphoneticized-source-field-name>
    <phoneticized-target-field-name>Person.LastName_Phon
  </phoneticized-target-field-name>
  <encoding-type>NYSIIS</encoding-type>
</phoneticize-field>
<phoneticize-field>
  <unphoneticized-source-field-name>Person.Alias[*].FirstName_Std
  </unphoneticized-source-field-name>
  <phoneticized-target-field-name>Person.FirstName_Phon
  </phoneticized-target-field-name>
  <encoding-type>Soundex</encoding-type>
</phoneticize-field>
<phoneticize-field>
  <unphoneticized-source-field-name>Person.Alias[*].LastName_Std
  </unphoneticized-source-field-name>
  <phoneticized-target-field-name>Person.LastName_Phon
  </phoneticized-target-field-name>
  <encoding-type>NYSIIS</encoding-type>
</phoneticize-field>
<phoneticize-field>
  <unphoneticized-source-field-name>
    Person.Address[*].AddressLine1_StName
  </unphoneticized-source-field-name>
  <phoneticized-target-field-name>
    Person.Address[*].AddressLine1_StPhon
  </phoneticized-target-field-name>
  <encoding-type>NYSIIS</encoding-type>
</phoneticize-field></phoneticize-fields>

```

## Configuring the Match String for Person Data (Repository)

You can include any fields on which you want to match in the match string. The match string is defined by the *match-column* elements in the MatchingConfig section of the Match Field file. If you specify a match type for a field in the wizard, that field (or any fields parsed from that field) is automatically defined in the match string.

To configure the match string, follow the instructions under “Defining the Master Index Match String (Repository)” in *Configuring Sun Master Indexes (Repository)*. For the Sun Match Engine, each data type has a different match type (specified by the *match-type* element). The FirstName and LastName match types are specific to person matching. You can specify any of the other match types defined in the match configuration file as well. For more information, see [“Sun Match Engine Match and Standardization Types” on page 24](#).

A sample match string for person matching is shown below. This sample matches on first and last names, date of birth, social security number, gender, and the street name of the address.

```

<match-system-object>
  <object-name>Person</object-name>

```

```
<match-columns>
  <match-column>
    <column-name>
      Enterprise.SystemSBR.Person.FirstName_Std
    </column-name>
    <match-type>FirstName</match-type>
  </match-column>
  <match-column>
    <column-name>Enterprise.SystemSBR.Person.LastName_Std
    </column-name>
    <match-type>LastName</match-type>
  </match-column>
  <match-column>
    <column-name>Enterprise.SystemSBR.Person.SSN
    </column-name>
    <match-type>SSN</match-type>
  </match-column>
  <match-column>
    <column-name>Enterprise.SystemSBR.Person.DOB
    </column-name>
    <match-type>DateDays</match-type>
  </match-column>
  <match-column>
    <column-name>Enterprise.SystemSBR.Person.Gender
    </column-name>
    <match-type>Char</match-type>
  </match-column>
  <match-column>
    <column-name>Enterprise.SystemSBR.Person.Address.StreetName
    </column-name>
    <match-type>StreetName</match-type>
  </match-column>
</match-columns>
</match-system-object>
```

## Sun Match Engine Address Data Type Configuration

Processing street addresses involves parsing, normalizing, and phonetically encoding certain fields prior to matching. The following topics describe the configuration files that define address processing logic and provide instructions for modifying the Match Field file for processing address fields.

- [“Sun Match Engine Address Matching Overview”](#) on page 49
- [“Match Configuration for Address Data \(Repository\)”](#) on page 51
- [“Sun Match Engine Standardization Configuration for Address Data”](#) on page 51
- [“Modifying Sun Match Engine Address Data Configuration Files”](#) on page 62



- [“Configuring the Matching Service for Address Data \(Repository\)” on page 63](#)

## Sun Match Engine Address Matching Overview

Matching on the address data type includes both standardizing and matching on address information in the master index application. You can implement street address standardization and matching on its own, or within a master index application designed to process person or business information. For example, standardizing address information allows you to include address fields as search criteria, even though matching might not be performed against these fields.

The Sun Match Engine can create standardized and phonetic values for street address information. Several configuration files are designed specifically to handle address data to define additional logic for the standardization and phonetic encoding process. These include address clues files, a patterns file, and a constants file. The United States address standardization engine is based on the work performed at the US Census Bureau. The clues files, in particular, are based on census bureau statistics.

The Sun Match Engine can match on any field as long as the match type for the field is defined in the match configuration file (`matchConfigFile.cfg`).

For information about the fields involved in address standardization and matching, see [“Sun Match Engine Address Data Processing Fields” on page 49](#).

## Sun Match Engine Address Data Processing Fields

When matching on address data, not all fields in a record need to be processed by the Sun Match Engine. The match engine only needs to process address fields that must be parsed, normalized, or phonetically encoded, and the fields against which matching is performed. These fields are defined in the Match Field file and processing logic for each field is defined in the standardization and matching configuration files.

### Address Data Match String Fields

The match string processed by the Sun Match Engine is defined by the match fields specified in the Match Field file. If you specify an “Address” match type for any field in the wizard, the default fields that store the parsed data are automatically added to the match string in the Match Field file. These fields include the house number, street direction, street type, and street name. You can remove any of these fields from the match string.

The match engine can process any combination of fields you specify for matching. By default, the match configuration file (`matchConfigFile.cfg`) includes rows specifically for matching on the fields that are parsed from the street address fields, such as the street number, street direction, and so on. The file also defines several generic match types. You can use any of the existing rows for matching or you can add rows for the fields you want to match.

## Address Data Standardized Fields

The Sun Match Engine expects that street address data will be provided in a free-form text field containing several components that must be parsed. By default, the match engine is configured to parse these components and to normalize and phonetically encode the street name. You can specify additional fields for phonetic encoding.

If you specify an “Address” match type for any field in the wizard, a standardization structure for that field is defined in the Match Field file. The fields listed below under “[Address Data Object Structure](#)” on page 50 are automatically defined as the target fields. Each of these fields has several entries in the standardization structure. This is because different parsed components can be stored in the same field. For example, the house number, post office box number, and rural route identifier are all stored in the house number field. If you do not specify address fields for matching in the wizard but want to standardize the fields, you can create a standardization structure in the Match Field file.

## Address Data Object Structure

The address fields specified for standardization are parsed into several additional fields. If you specify the “Address” match type in the wizard, the following fields are automatically added to the object structure and database creation script.

- *field\_name\_HouseNo*
- *field\_name\_StName*
- *field\_name\_StDir*
- *field\_name\_StType*
- *field\_name\_StPhon*

where *field\_name* is the name of the field for which you specified address matching. For example, if you specify the Address match type for the AddressLine1 field, the following fields are automatically added to the structure: AddressLine1\_HouseNo, AddressLine1\_StName, AddressLine1\_StDir, AddressLine1\_StType, and AddressLine1\_StPhon.

You can add these fields manually if you do not specify a match type in the wizard.

---

**Note** – The object structure for Sun Master Patient Index uses a slightly different naming convention. For the names of the fields defined for Sun Master Patient Index, refer to the *Sun SeeBeyond eIndex Single Patient View User’s Guide*.

---

## Match Configuration for Address Data (Repository)

The default match configuration file, `matchConfigFile.cfg`, defines several match types for the kinds of address data typically included in the match string. You can customize the existing match types or create new match types for the data being processed. The following match types are typical for matching on address data.

- 
- |                |              |
|----------------|--------------|
| ▪ StreetNumber | ▪ StreetDir  |
| ▪ HouseNumber  | ▪ StreetType |
- 

In addition, you can use any of these generic match types for matching on address data.

- 
- |           |        |
|-----------|--------|
| ▪ String  | ▪ SSN  |
| ▪ Date    | ▪ Char |
| ▪ Numeric | ▪ pro  |
| ▪ Integer | ▪ Exac |
| ▪ Real    |        |
- 

The match configuration file appears under the Match Engine node of the master index project. For more information about the comparison functions used for each match type and how the weights are tuned, see [“Customizing the Match Configuration” on page 88](#) and [“Match Configuration Comparison Functions for Sun Match Engine \(Repository\)” on page 93](#).

## Sun Match Engine Standardization Configuration for Address Data

Several configuration files are used to define address processing logic for the Sun Match Engine. You can customize any of the configuration files described in this section to fit your processing and standardization requirements for address data. There are no address standardization files that are common to all domains; all address files are domain-specific. These files are located within the domain-specific folders of the Standardization Engine node (with two exceptions noted below).

Address standardization files are specific to each domain and include patterns and clues files, as well as files that define internal and external constants. The domain corresponding to each file is indicated at the end of the file name; for example, `addressConstantsUK.cfg` and `addressConstantsFR.cfg`. These domain abbreviations are indicated by an asterisk (\*) in the file descriptions in the following topics.

- “The Address Constants File (addressConstants\*.cfg)” on page 52
- “The Address Clues File (addressClueAbbrev\*.dat)” on page 53
- “The Address Internal Constants File (addressInternalConstants\*.cfg)” on page 54
- “The Address Master Clues File (addressMasterClues\*.dat)” on page 54
- “The Address Patterns File (addressPatterns\*.dat)” on page 56
- “The Address Output Patterns File (addressOutPatterns\*.dat)” on page 57
- “Address Pattern File Components” on page 59

## The Address Constants File (addressConstants\*.cfg)

The address constants file defines certain information about the standardization files used for processing address data, primarily the number of lines contained in each file. The number of lines specified here must be equal to or greater than the number of lines actually contained in each file. The constants file for United States data is in the Standardization node of the project and is named `addressConstants.cfg`; the constants file for the other domains is located under the domain name node.

Table 15 lists and describes each parameter in the constants file. The files referenced by these parameters are described on the following pages.

TABLE 15 Address Constants File Parameters

Parameter	Description
<code>maxWords</code>	The maximum number of words in a given address field.
<code>clueArraySize</code>	The maximum number of lines in the address clues file ( <code>addressClueAbbrev*.dat</code> ).
<code>patternArraySize</code>	The maximum number of lines in the patterns file ( <code>addressPatterns*.dat</code> ).
<code>maxPattSize</code>	The maximum length (in characters) of any pattern in the address patterns file.
<code>imageSize</code>	The maximum length of an input address field.
<code>nameOutputFieldSize</code>	The maximum output length of a street or property name.
<code>numberOutputFieldSize</code>	The maximum output length of a house number or rural route number within the structure identifier or post office box fields.
<code>directionOutputFieldSize</code>	The maximum output length of a directional field (prefix or suffix).
<code>typeOutputFieldSize</code>	The maximum output length of a street type field (prefix or suffix).
<code>prefixOutputFieldSize</code>	The maximum length of a number prefix field.
<code>suffixOutputFieldSize</code>	The maximum length of a number suffix field.
<code>extensionOutputFieldSize</code>	The maximum output length of any extension field.

TABLE 15 Address Constants File Parameters (Continued)

Parameter	Description
extrainfoOutputFieldSize	The maximum output length of any miscellaneous information that is not recognized as a known type.

## The Address Clues File (addressClueAbbrev\*.dat)

The address clues file lists common terms in street addresses, specifies a normalized value for each common term, and categorizes the terms into street address component types. A term can be categorized into multiple component types. The relevance value specifies which of the component types the term is most likely to be. For example, the term “Junction” is standardized as “Jct”, and is classified as a street type, building unit, and generic term (giving relevance in that order).

This file helps the Sun Match Engine recognize common terms in street addresses, and to parse and normalize the values correctly. The syntax of this file is:

common-term normalized-term ID-number/type-token

You can modify or add entries in this table as needed. [Table 16](#) describes the columns in the addressClueAbbrev\*.dat file.

TABLE 16 Address Clues File Columns

Column	Description
common-term	A term commonly found in street addresses.
normalized-term	The normalized version of the common term.
ID-number/type-token	An ID number and a token indicating the type of address component represented by the common term. The ID number corresponds to an ID number in the address master clues file, and the type token corresponds to the type specified for that ID number in the address master clues file. One term might have several ID number and token type pairs.

Following is an excerpt from the addressClueAbbrevUS.dat file.

TRLR VLG	Trpk	59BU		
TRPK	Trpk	59BU		
TRPRK	Trpk	59BU		
VILLA	Vlla	305TY	60BU	
VLLA	Vlla	305TY	60BU	
VILLAS	Vlla	60BU		
VILL	Vlg	317TY	61BU	364AU
VILLAG	Vlg	317TY	61BU	364AU
VLG	Vlg	317TY	61BU	364AU
VILLAGE	Vlg	317TY	61BU	364AU

VILLG	Vlg	317TY	61BU	364AU
VILLIAGE	Vlg	317TY	61BU	364AU
VLGE	Vlg	317TY	61BU	364AU
VIVI	Vivi	62BU		
VIVIENDA	Vivi	62BU		
COLLEGE	Coll	64BU		0AU
CLG	Coll	64BU		
COTTAGE	Cott	65BU	65BP	0AU

## The Address Internal Constants File (`addressInternalConstants*.cfg`)

The address internal constants file defines and configures tokens and array sizes used by the address standardizer. This file is used internally by the standardization engine and most of the parameters should not be modified.

One parameter you might need to modify is *spCh*, which defines any special characters that should not be removed from addresses during standardization. By default, the standardization process keeps hyphens (-), pound signs (#), forward slashes (/), ampersands (&), and pipes (|). Any other special characters found in the address are removed unless they are defined for the *spCh* parameter. Delineate each special character in the list with a space, as shown below.

```
spCh = & < >
```

Characters that are not included in the standard ISO 8859-1 (Latin-1) character set must be preceded by a back slash (\) and represented in Unicode. For example, use the following to retain right and left single quotes (” ’) in addresses:

```
spCh = \u2018 \u2019
```

---

**Note** – Periods (.) and commas (,) are always removed from addresses, even if they are added to the *spCh* list.

---

## The Address Master Clues File (`addressMasterClues*.dat`)

The address master clues file lists common terms in street addresses as defined by the United States Postal Service (USPS), the United Kingdom’s Royal Mail, the Australian Postal Corporation, or France’s La Poste (depending on the domain in use). For each common term, this file specifies a normalized value, defines postal information, and categorizes the terms into street address component types. A term can be categorized into multiple component types.

The syntax of this file is:

```
ID-number common-term normalized-term short-abbrev postal-abbrev CFCCS
type-token usage-flag postal-flag
```

You can modify or add entries in this table as needed. [Table 17](#) describes the columns in the `addressMasterClues*.dat` file.

TABLE 17 Address Master Clue File Columns

Column	Description
ID-number	A unique identification number for the address common term. This number corresponds to an ID number for the same term in the address clues file.
common-term	A common address term, such as Park, Village, North, Route, Centre, and so on.
normalized-term	The normalized version of the common term.
short-abbrev	A short abbreviation of the common term.
postal-abbrev	The standard postal abbreviation of the common term.
CFCCS	The census feature class code of the term (as defined in the Census Tiger® database). The following values are used: <ul style="list-style-type: none"> <li>■ A – Road</li> <li>■ B – Railroad</li> <li>■ C – Miscellaneous</li> <li>■ D – Landmark</li> <li>■ E – Physical feature</li> <li>■ F – Nonvisible feature</li> <li>■ H – Hydrography</li> <li>■ X – Unclassified</li> </ul>
type-token	The type of address component represented by the common term. Types are specified by an address token (for more information, see <a href="#">“Address Type Tokens” on page 59</a> ).
usage-flag	A flag indicating how the term is used (for more information, see <a href="#">“Pattern Classes” on page 61</a> ).
postal-flag	The standard postal code for the term.

Following is an excerpt from the addressMasterCluesUS.dat file.

11Alley	Alley	Al	Aly A	TY R U
12Alternate Route	Alt Rte	Alt	Alt A	TY R
15Arcade	Arcade	Arc	Arc A	TY R U
16Arroyo	Arroyo	Arryo	ArryHA	TY R
17Autopista	Atpta	Apta	AptaA	TY R
18Avenida	Avenida	Ava	Ava A	TY R
19Avenue	Avenue	Ave	Ave A	TY R U
26Boulevard	Blvd	Blvd	BlvdA	TY R U
32Bulevar	Blvr	Blv	Blv A	TY R
33Business Route	Bus Rte	BusRt	BsRtA	TY R
34Bypass	Bypass	Byp	Byp A	TY R U
36Calle	Calle	Calle	ClleA	TY R
37Calleja	Calleja	Cja	Cja A	TY R

38Callejon	Callej	Cjon	CjonA	TY R
39Camino	Camino	Cam	Cam A	TY R
47Carretera	Carrrt	Carr	CarrA	TY R
48Causeway	Cswy	Cswy	CswyAH	TY R U
51Center	Center	Ctr	Ctr DA	TY R U

## The Address Patterns File (addressPatterns\*.dat)

The address patterns file defines the expected input patterns of each individual street address field being standardized so the Sun Match Engine can recognize and process these values.

Tokens are used to indicate the type of address component in the input and output fields. This file contains two rows for each pattern. The first row defines the input pattern for each address field and provides an example. The second row defines the output pattern for each address field, the pattern type, the relative importance of the pattern compared to other patterns, and usage flags (as shown below).

```
AU A1 TY           01 Oak B Street
NA NA ST           T* 75             TX
```

When an address is parsed, each line of the address is delineated by a pipe (|) and sent to the parser separately. The output tokens for each line are then concatenated and the output pattern is processed using the addressOutPatterns\*.dat file to determine whether the output pattern is listed in the file. If the pattern is found, output patterns are modified as indicated in the addressOutPatterns\*.dat file to resolve any ambiguities that might arise when two lines of address information contain common elements. The relative importance determines which pattern to use when the format of the input field matches more than one pattern. This file should only be modified by personnel with a thorough understanding of address patterns and tokens.

The syntax of this file is:

```
input-pattern example output-pattern pattern-class pattern-modifier priority
usage-flag exclude-flag
```

You can modify or add entries in this table as needed. [Table 18](#) describes the columns in the addressPatterns\*.dat file.

**TABLE 18** Address Patterns File

Column	Description
input-pattern	Tokens that represent a possible input pattern from an individual unparsed street address field. Each token represents one component. For more information about address tokens, see <a href="#">“Address Type Tokens” on page 59</a> .



TABLE 18 Address Patterns File (Continued)

Column	Description
example	An example of a street address that fits the specified pattern. This file element is optional.
output-pattern	Tokens that represent the output pattern for the specified input pattern. Each token represents one component of the output of the Sun Match Engine. For more information about address tokens, see <a href="#">“Address Type Tokens” on page 59.</a>
pattern-class	An indicator of the type of address component represented by the pattern. Possible pattern types are listed in <a href="#">“Pattern Classes” on page 61</a> <a href="#">“Pattern Classes” on page 61.</a>
pattern-modifier	An indicator of whether the priority of the pattern is averaged against other patterns that match the input. Pattern modifiers are listed in <a href="#">“Pattern Modifiers” on page 62.</a>
priority	The priority weight to use for the pattern when the pattern is a sub-pattern of a larger input pattern. For more information, see <a href="#">“Priority Indicators” on page 62.</a>
usage-flag	A flag indicating how the term is used (for more information, see <a href="#">“Pattern Classes” on page 61</a> ). This file element is optional.
exclude-flag	This file element is optional.

Following is an excerpt from the addressPatternsUS.dat file.

```

NU DR TY A1 AU           01  123 South Avenida B Oak
HN PD PT NA NA          H* 70

NU DR TY NU DR           01  123 South Avenida 1 West
HN PD PT NA SD          H* 70

NU A1 TY AU TY           01  123 C circle hill drive
HN HS NA NA ST          H* 70

NU A1 AM A1 TY           01  123 M & M road
HN NA NA NA ST          H* 65

NU TY AU A1              01  123 Avenida Oak B
HN PT NA NA              H* 60

NU TY NU A1              01  123 Avenida 1 B
HN PT NA NA              H* 60

```

## The Address Output Patterns File (addressOutPatterns\*.dat)

The address output patterns file uses the field patterns output by the addressPatterns\*.dat file to determine how to parse all standardized address fields. As with the

addressPatterns\*.dat file, tokens are used to indicate the type of address component in the input and output data. This file contains two rows for each pattern. The first row defines the input pattern received from addressPatterns\*.dat and provides an example. The second row defines the output pattern (as shown below).

```
EI|BN BT|*           // HILLVIEW|FULBOURN HOSPITAL
BN|BI BY
```

The syntax of this file is:

```
input-pattern        example
output-pattern
```

You can modify or add entries in this table as needed. [Table 19](#) describes the columns in the addressOutputPatterns\*.dat file.

**TABLE 19** Address Output Patterns File

Column	Description
input-pattern	Tokens that represent a possible input pattern from addressPatterns*.dat. Each token represents one component and the pattern for each address field in the address is separated by a pipe ( ). For more information about address tokens, see “Address Type Tokens” on page 59. Note that this file only uses output tokens.
example	An example of a street address that fits the specified pattern. This file element is optional.
output-pattern	Tokens that represent the output pattern for the specified input pattern. Each token represents one component of the output of the Sun Match Engine. For more information about address tokens, see “Address Type Tokens” on page 59.

Following is an excerpt from the addressPatternsUS.dat file. In the first example, addressPatternsUS.dat outputs three address fields containing these components: building name and type; street name and type; and street name and type. addressOutputPatternsUS.dat changes the tokens for the second street name and type to indicate they are not the primary street name and type. Therefore, “New Bridge” is populated into the parsed street name field in the database.

```
BN BT|NA ST|NA ST|*           // PROTEA HOUSE|NEW BRIDGE|MARINE PARADE
BN BT|NA ST|N2 S2

HN NA ST|HN NA ST|*           // 21 HEIGHWAY COURT|45 BROOKLAND ROAD
HN NA ST|H2 N2 S2

HN NA ST|NA ST|*             // 21 HEIGHWAY COURT|BROOKLAND ROAD
HN NA ST|N2 S2
```

```
NA ST|HN NA ST|*           // HEIGHWAY COURT|45 BROOKLAND ROAD
NA ST|H2 N2 S2
```

## Address Pattern File Components

The address patterns files use pattern type tokens, pattern classes, pattern modifies, and priority indicators to process and parse address data. Before modifying any of the patterns files, you must have a good understanding of these file components.

## Address Type Tokens

The address pattern and clues files use tokens to denote different components in a street address, such as street type, house number, street names, and so on. These files use one set of tokens for input fields and another set for output fields. You can use only the predefined tokens to represent address components; the Sun Match Engine does not recognize custom tokens.

[Table 20](#) lists and describes each input token; [Table 21](#) lists and describes each output token.

**TABLE 20** Input Address Pattern Type Tokens

Token	Description
A1	Alphabetic value, one character in length
AM	Ampersand
AU	Generic word
BP	Building property
BU	Building unit
BX	Post office box
DA	Dash (as a starting character)
DR	Street direction
EI	Extra information
EX	Extension
FC	Numeric fraction
HR	Highway route
MP	Mile posts
NL	Common words, such as “of”, “the”, and so on
NU	Numeric value

TABLE 20 Input Address Pattern Type Tokens (Continued)

Token	Description
OT	Ordinal type
PT	Prefix type
RR	Rural route
SA	State abbreviation
TY	Street type
WD	Descriptor within the structure
WI	Identifier within the structure

Table 21 lists and describes each output token.

TABLE 21 Output Address Pattern Tokens

Token	Description
1P	Building number prefix
2P	Second building number prefix
BD	Property or building directional suffix
BI	Structure (building) identifier
BN	Property or building name
BS	Building number suffix
BT	Property or building type suffix
BX	Post office box descriptor
BY	Structure (building) descriptor
DB	Property or building directional prefix
EI	Extra information
EX	Extension index
H1	First house number (the actual number)
H2	Second house number (house number suffix)
HN	House number
HS	House number suffix

TABLE 21 Output Address Pattern Tokens (Continued)

Token	Description
N2	Second street name
NA	Street name
NB	Building number
NL	Conjunctions that connect words or phrases in one component type (usually the street name)
P1	House number prefix
P2	Second house number prefix
PD	Directional prefix to the street name
PT	Street type prefix to the street name
RR	Rural route descriptor
RN	Rural route identifier
S2	Street type suffix to the second street name
SD	Directional suffix to the street name
ST	Street type suffix to the street name
TB	Property or building type prefix
WI	Identifier within the structure
WD	Descriptor within the structure
XN	Post office box identifier

## Pattern Classes

Each pattern defined in the address patterns file must have an associated pattern class. The pattern class indicates a portion of the input pattern or the type of address data that is represented by the pattern. You can specify any of the following pattern classes.

- **H** - the address pattern represents a house
- **B** - the address pattern represents a building
- **W** - the address pattern represents a unit within a structure, such as an apartment or suite number
- **T** - the address pattern represents a street type or direction
- **R** - the address pattern represents a rural route
- **P** - the address pattern represents a Post Office box
- **N** - the address pattern is mostly numeric

These classes are also specified as usage flags in the patterns file and the master clues file.

## Pattern Modifiers

Each pattern type must be followed by a pattern modifier that indicates how to handle cases where one or more defined patterns is found to be a sub-pattern of a larger input pattern. In this case, the Sun Match Engine must know how to prioritize each defined pattern that is a part of the larger pattern. There are two pattern modifiers.

- \* - An asterisk indicates that the priority weight for the matching pattern is averaged down equally with the other matching sub-patterns.
- + - A plus sign indicates that the priority weight for the matching pattern is not averaged down equally with the other matching sub-patterns.

## Priority Indicators

The priority indicator is a numeric value following the pattern modifier that indicates the priority weight of the pattern. These values work best when defined as a multiple of five between and including 35 and 95. If a pattern is assigned a priority of 90 or 95 and the pattern matches, or is a sub-pattern of, the input pattern, the match engine stops searching for additional matching patterns and uses the high-priority matching pattern.

# Modifying Sun Match Engine Address Data Configuration Files

To customize the Sun Match Engine configuration files for processing street address data, you can modify any of the address standardization files using the text editor provided in NetBeans. Before modifying the match configuration file, review the information provided in [“Sun Match Engine Matching Configuration” on page 12](#) and [“Match Configuration Comparison Functions for Sun Match Engine \(Repository\)” on page 93](#). Make sure a thorough data analysis has been performed to determine the best fields for matching and the best comparison functions to use for each field.

Updating most standardization files is a straightforward process. Make sure to follow the syntax guidelines provided in [“Sun Match Engine Standardization Configuration for Address Data” on page 51](#). If you add rows to any of the standardization files, make sure to adjust the corresponding parameter in the address constants file (`addressConstants.cfg`).

Modifying the patterns file is a more complex task. Only modify this file once you fully understand pattern tokens, types, relevance, and flags.

## Configuring the Matching Service for Address Data (Repository)

To ensure the master index application uses the Sun Match Engine to process address information, you must customize the Matching Service. This includes modifying the Match Field file to support the fields on which you want to match, to standardize the appropriate fields, and to specify the Sun Match Engine as the match and standardization engine (by default, the Sun Match Engine is already specified so this does not need to be changed). Perform the following tasks to configure the Matching Service.

- “Configuring the Standardization Structure for Address Data (Repository)” on page 63
- “Configuring the Match String for Address Data (Repository)” on page 66

When configuring the Matching Service, keep in mind the information presented in “Configuring the Master Index Matching Service (Repository)” on page 27.

## Configuring the Standardization Structure for Address Data (Repository)

The standardization structure is configured in the StandardizationConfig section of the Match Field file, which is described in detail in “Match Field Configuration (Repository)” in *Understanding Sun Master Index Configuration Options (Repository)*. To configure the required fields for standardization and phonetic encoding, modify the standardization and phonetic encoding structures. The following sections provide additional guidelines and samples specific to standardizing address data.

---

**Note** – In the default configuration, the rules defined for the address data type assume that all input fields must be parsed as well as normalized. Thus, there is no need to configure fields only for normalization.

---

### Address Standardization Structures

For address fields, the source fields in the standardization structure must include the fields predefined for parsing and normalization. This includes any fields containing street address information, which are parsed into the street address fields listed in “Address Data Object Structure” on page 50 (excluding the phonetic street name field). The target fields can include any of these parsed fields. Follow the instructions under “Defining Master Index Standardization Rules (Repository)” in *Configuring Sun Master Indexes (Repository)* to define fields for normalization. For the *standardization-type* element, enter **Address** (for more information, see “Sun Match Engine Match and Standardization Types” on page 24). For a list of field IDs to use in the *standardized-object-field-id* element, see Table 3.

A sample standardization structure for address data is shown below. This structure parses the first two lines of street address into the standard street address fields. Only the United States domain is defined in this structure.

```
free-form-texts-to-standardize>
<group standardization-type="ADDRESS"
  domain-selector="com.stc.eindex.matching.impl.SingleDomainSelectorUS">
  <unstandardized-source-fields>
    <unstandardized-source-field-name>Person.Address[*].Address1
  </unstandardized-source-field-name>
    <unstandardized-source-field-name>Person.Address[*].Address2
  </unstandardized-source-field-name>
  </unstandardized-source-fields>
  <standardization-targets>
    <target-mapping>
      <standardized-object-field-id>HouseNumber
    </standardized-object-field-id>
      <standardized-target-field-name>Person.Address[*].HouseNumber
    </standardized-target-field-name>
    </target-mapping>
    <target-mapping>
      <standardized-object-field-id>RuralRouteIdentif
    </standardized-object-field-id>
      <standardized-target-field-name>Person.Address[*].HouseNumber
    </standardized-target-field-name>
    </target-mapping>
    <target-mapping>
      <standardized-object-field-id>BoxIdentif
    </standardized-object-field-id>
      <standardized-target-field-name>Person.Address[*].HouseNumber
    </standardized-target-field-name>
    </target-mapping>
    <target-mapping>
      <standardized-object-field-id>MatchStreetName
    </standardized-object-field-id>
      <standardized-target-field-name>Person.Address[*].StreetName
    </standardized-target-field-name>
    </target-mapping>
    <target-mapping>
      <standardized-object-field-id>RuralRouteDescript
    </standardized-object-field-id>
      <standardized-target-field-name>Person.Address[*].StreetName
    </standardized-target-field-name>
    </target-mapping>
    <target-mapping>
      <standardized-object-field-id>BoxDescript
    </standardized-object-field-id>
      <standardized-target-field-name>Person.Address[*].StreetName
    </standardized-target-field-name>
    </target-mapping>
    <target-mapping>
      <standardized-object-field-id>PropDesPrefDirection
```



```

        </standardized-object-field-id>
        <standardized-target-field-name>Person.Address[*].StreetDir
        </standardized-target-field-name>
    </target-mapping>
    <target-mapping>
        <standardized-object-field-id>PropDesSufDirection
        </standardized-object-field-id>
        <standardized-target-field-name>Person.Address[*].StreetDir
        </standardized-target-field-name>
    </target-mapping>
    <target-mapping>
        <standardized-object-field-id>StreetNameSufType
        </standardized-object-field-id>
        <standardized-target-field-name>Person.Address[*].StreetType
        </standardized-target-field-name>
    </target-mapping>
    <target-mapping>
        <standardized-object-field-id>StreetNamePrefType
        </standardized-object-field-id>
        <standardized-target-field-name>Person.Address[*].StreetType
        </standardized-target-field-name>
    </target-mapping>
</standardization-targets>
</group>
</free-form-texts-to-standardize>

```

## Address Phonetic Encoding

When you match or standardize on street address fields, the street name should be specified for phonetic conversion (this is done by default). Follow the instructions under “Defining Phonetic Encoding for the Master Index (Repository)” in *Configuring Sun Master Indexes (Repository)* to define fields for phonetic encoding.

A sample of the *phoneticize-fields* element is shown below. This sample only converts the address street name. You can define additional fields for phonetic encoding.

```

<phoneticize-fields>
    <phoneticize-field>
        <unphoneticized-source-field-name>Person.Address[*].StreetName
        </unphoneticized-source-field-name>
        <phoneticized-target-field-name>Person.Address[*].StreetName_Phon
        </phoneticized-target-field-name>
        <encoding-type>NYSIIS</encoding-type>
    </phoneticize-field>
</phoneticize-fields>

```

## Configuring the Match String for Address Data (Repository)

For matching on street address fields, make sure the match string you specify in the MatchingConfig section of the Match Field file contains all or a subset of the fields that contain the standardized data (the original text in street address fields are generally too inconsistent to use for matching). You can include additional fields for matching, such as the city name or postal code.

To configure the match string, follow the instructions under “Defining the Master Index Match String (Repository)” in *Configuring Sun Master Indexes (Repository)*. For the Sun Match Engine, each component of a street address has a different match type (specified by the *match-type* element). The default match types for addresses are StreetName, HouseNumber, StreetDir, and StreetType. You can specify any of the other match types defined in the match configuration file, as well. For more information, see [“Sun Match Engine Match and Standardization Types” on page 24](#).

A sample match string for address matching is shown below.

```
<match-system-object>
  <object-name>Person</object-name>
  <match-columns>
    <match-column>
      <column-name>Enterprise.SystemSBR.Person.Address.StreetName
      </column-name>
      <match-type>StreetName</match-type>
    </match-column>
    <match-column>
      <column-name>Enterprise.SystemSBR.Person.Address.HouseNumber
      </column-name>
      <match-type>HouseNumber</match-type>
    </match-column>
    <match-column>
      <column-name>Enterprise.SystemSBR.Person.Address.StreetDir
      </column-name>
      <match-type>StreetDir</match-type>
    </match-column>
    <match-column>
      <column-name>Enterprise.SystemSBR.Person.Address.StreetType
      </column-name>
      <match-type>StreetType</match-type>
    </match-column>
  </match-columns>
</match-system-object>
```

# Sun Match Engine Business Names Data Type Configuration

Processing business name fields involves parsing, normalizing, and phonetically encoding certain fields prior to matching. The following topics describe the configuration files that define business name processing logic and provide instructions for modifying the Match Field file for processing business names.

- [“Sun Match Engine Business Name Matching Overview”](#) on page 67
- [“Sun Match Engine Match Configuration for Business Names”](#) on page 69
- [“Sun Match Engine Standardization Configuration for Business Names”](#) on page 69
- [“Modifying Sun Match Engine Business Name Configuration Files”](#) on page 83
- [“Configuring the Matching Service for Business Names \(Repository\)”](#) on page 83

## Sun Match Engine Business Name Matching Overview

Matching on the business name data type includes standardizing and matching on free-form business name fields. You can implement business name standardization and matching on its own or within a master index application designed to process person information. For example, standardizing business name fields allows you to include these fields as search criteria, even though matching might not be performed against these fields.

The Sun Match Engine can create standardized and phonetic values for business names. Several configuration files are designed specifically to handle business names to define additional logic for the standardization and phonetic encoding process. These include reference files, a patterns file, and key type files. The Sun Match Engine can match on any field as long as the match type for the field is defined in the match configuration file (`matchConfigFile.cfg`). The business name standardization files are common to all national domains, so no domain-specific configuration is required.

For more information about the fields involved in business name standardization and matching, see [“Sun Match Engine Business Name Processing Fields”](#) on page 67.

## Sun Match Engine Business Name Processing Fields

When matching on free-form business names, not all fields in a record need to be processed by the Sun Match Engine. The match engine only needs to process fields that must be parsed, normalized, or phonetically converted, and the fields against which matching is performed. These fields are defined in the Match Field file, and processing logic for each field is defined in the standardization and matching configuration files.

## Business Name Match String Fields

The match string processed by the Sun Match Engine is defined by the match fields specified in the Match Field file. If you specify a “BusinessName” match type for any field in the wizard,

most of the parsed business name fields are automatically added to the match string in the Match Field file, including the name, organization type, association type, sector, industry, and URL. You can remove any of these fields from the match string.

The match engine can process any combination of fields you specify for matching. By default, the match configuration file (`matchConfigFile.cfg`) includes rows specifically for matching on the fields that are parsed from the business name fields. The file also defines several generic match types. You can use any of the existing rows for matching or you can add rows for the fields you want to match.

## Business Name Standardized Fields

The Sun Match Engine expects that business name data will be provided in a free-form text field containing several components that must be parsed. The match engine is designed to parse these components, and to normalize and phonetically encode the business name. You can specify additional fields for phonetic encoding.

If you specify the “BusinessName” match type for any field in the wizard, a standardization structure for that field is defined in the Match Field file. The fields defined as the target fields are listed in the next section, “[Business Name Object Structure](#)” on page 68.

## Business Name Object Structure

For the default configuration of the business name data type, the address fields specified for standardization are parsed into several additional fields, one of which is also normalized. If you specify the appropriate match type in the wizard, the following fields are automatically added to the object structure and database creation script.

- *field\_name\_Name*
- *field\_name\_NamePhon*
- *field\_name\_OrgType*
- *field\_name\_AssocType*
- *field\_name\_Industry*
- *field\_name\_Sector*
- *field\_name\_Alias*
- *field\_name\_Url*

where *field\_name* is the name of the field for which you specified business name matching. For example, if you specify the BusinessName match type for the Company field, the fields automatically added to the structure include `Company_Name`, `Company_NamePhon`, `Company_OrgType`, and so on.

You can add these fields manually if you do not specify a match type in the wizard.

## Sun Match Engine Match Configuration for Business Names

The default match configuration file, `matchConfigFile.cfg`, defines several match types for the kinds of business name data typically included in the match string. You can customize the existing match types or create new match types for the data being processed. The following match types are typical for matching on business names.

- 
- |                      |                       |
|----------------------|-----------------------|
| ▪ PrimaryName        | ▪ AliasList           |
| ▪ OrgTypeKeyword     | ▪ IndustryTypeKeyword |
| ▪ AssocTypeKeyword   | ▪ URL                 |
| ▪ IndustrySectorList |                       |
- 

In addition, you can use any of these generic match types for matching on business names.

- 
- |           |        |
|-----------|--------|
| ▪ String  | ▪ Real |
| ▪ Date    | ▪ Char |
| ▪ Numeric | ▪ pro  |
| ▪ Integer | ▪ Exac |
- 

This file appears under the Match Engine node of the master index project. For more information about the comparison functions used for each match type and how the weights are tuned, see [“Customizing the Match Configuration” on page 88](#) and [“Match Configuration Comparison Functions for Sun Match Engine \(Repository\)” on page 93](#).

## Sun Match Engine Standardization Configuration for Business Names

Several configuration files are used to define business name processing logic for the Sun Match Engine. You can customize any of the configuration files described in this section to fit your data processing and standardization requirements. These files appear under the Standardization Engine node of the master index project.

The following topics described each file used for business name standardization:

- [“The Business Constants File \(bizConstants.cfg\)” on page 70](#)
- [“The Adjectives Key Type File \(bizAdjectivesTypeKeys.dat\)” on page 71](#)

- “The Alias Key Type File (bizAliasTypeKeys.dat)” on page 71
- “The Association Key Type File (bizAssociationTypeKeys.dat)” on page 72
- “The General Terms Reference File (bizBusinessGeneralTerms.dat)” on page 73
- “The City or State Key Type File (bizCityorStateTypeKeys.dat)” on page 73
- “The Business Former Name Reference File (bizCompanyFormerNames.dat)” on page 74
- “The Merged Business Name Category File (bizCompanyMergerNames.dat)” on page 74
- “The Primary Business Name Reference File (bizCompanyPrimaryNames.dat)” on page 75
- “The Connector Tokens Reference File (bizConnectorTokens.dat)” on page 76
- “The Country Key Type File (bizCountryTypeKeys.dat)” on page 76
- “The Industry Sector Reference File (bizIndustryCategoryCode.dat)” on page 77
- “The Industry Key Type File (bizIndustryTypeKeys.dat)” on page 78
- “The Organization Key Type File (bizOrganizationTypeKeys.dat)” on page 79
- “The Business Patterns File (bizPatterns.dat)” on page 79
- “The Special Characters Reference File (bizRemoveSpecChars.dat)” on page 82

## The Business Constants File (bizConstants.cfg)

The business constants file defines certain information about the standardization files used for processing business data, primarily the number of lines contained in each file. The number of lines specified must be equal to or greater than the number of lines actually contained in each file.

Table 22 lists and describes each parameter in the constants file. The files referenced by these parameters are described on the following pages.

TABLE 22 Business Constants File Parameters

Parameter	Description
cityMax	The maximum number of lines in the city or state key type file (bizCityorStateTypeKey.dat).
primaryMax	The maximum number of lines in the primary business names reference file (bizCompanyPrimaryNames.dat).
countryMax	The maximum number of lines in the country key type file (bizCountryTypeKeys.dat).
industryMax	The maximum number of lines in the industry key type file (bizIndustryTypeKeys.dat).
patternMax	The maximum number of lines in the business patterns file (bizPatterns.dat).
mergerMax	The maximum number of lines in the merged business name category file (bizCompanyMergerNames.dat).
adjectiveMax	The maximum number of lines in the adjective key type file (bizAdjectiveTypeKeys.dat).

TABLE 22 Business Constants File Parameters (Continued)

Parameter	Description
orgMax	The maximum number of lines in the organization key type file (bizOrganizationTypeKeys.dat).
assocMax	The maximum number of lines in the association key type file (bizAssociationTypeKeys.dat).
genTermMax	The maximum number of lines in the general terms reference file (bizBusinessGeneralTerms.dat).
charsMax	The maximum number of lines in the special characters reference file (bizRemoveSpecChars.dat).
bizMaxWords	The maximum number of tokens allowed in the input business name. If no value is defined for this parameter, the default is the value set for the <i>words</i> parameter in the personConstants.cfg file.

## The Adjectives Key Type File (bizAdjectivesTypeKeys.dat)

The adjectives key type file defines adjectives commonly found in business names so the Sun Match Engine can recognize and process these values as a part of the business name. This file contains one column with a list of commonly used adjectives, such as General, Financial, Central, and so on.

You can modify or add entries in this file as needed. Following is an excerpt from the bizAdjectivesTypeKeys.dat file.

```
DIGITAL
DIRECTED
DIVERSIFIED
EDUCATIONAL
ELECTROCHEMICAL
ENGINEERED
EVOLUTIONARY
EXTENDED
FACTUAL
FEDERAL
```

## The Alias Key Type File (bizAliasTypeKeys.dat)

The alias key type file lists business name acronyms and abbreviations along with their standardized names so the Sun Match Engine can recognize and process these values correctly. You can add entries to the alias key type file using the following syntax.

```
alias standardized-name
```

Table 23 describes the columns in the bizAliasTypeKeys.dat file.

TABLE 23 Alias Key Type File

Column	Description
alias	An abbreviation or acronym commonly used in place of a specific business name.
standardized-name	The normalized version of the alias name.

Following is an excerpt from the `bizAliasTypeKeys.dat` file.

```

BBH          BARTLE BOGLE HEGARTY
BBH          BROWN BROTHERS HARRIMAN
IBM          INTERNATIONAL BUSINESS MACHINE
IDS          INCOMES DATA SERVICES
IDS          INSURANCE DATA SERVICES
IDS          THE INTEGRATED DECISION SUPPORT GROUP
IDS          THE INTERNET DATABASE SERVICE
CAL-TECH    CALIFORNIA INSTITUTE OF TECHNOLOGY

```

## The Association Key Type File (`bizAssociationTypeKeys.dat`)

The association key type file lists business association types along with their standardized names so the Sun Match Engine can recognize and process these values correctly. You can add entries to the association key type file using the following syntax.

```
association-type standardized-type
```

[Table 24](#) describes the columns in the `bizAssociationTypeKeys.dat` file.

TABLE 24 Association Type Key Table

Column	Description
association-type	A common association type for businesses, such as Partners, Group, and so on.
standardized-type	The standardized version of the association type. If this column contains a name instead of a zero, that name must also be listed in a different entry as an association type with a standardized form of "0".

Following is an excerpt from the `bizAssociationTypeKeys.dat` file.

```

ASSOCIATES    0
BANCORP       0
BANCORPORATION BANCORP
COMPANIES     0
GP            GROUP
GROUP         0
PARTNERS      0

```



## The General Terms Reference File (`bizBusinessGeneralTerms.dat`)

The general terms reference file lists terms commonly used in business names. This file is used to identify terms that indicate a business, such as bank, supply, factory, and so on, so the Sun Match Engine can recognize and process the business name.

This file contains one column that lists common terms in the business names you process. You can add entries as needed. Below is an excerpt from the `bizBusinessGeneralTerms.dat` file.

```
BUILDING
CITY
CONSUMER
EAST
EYE
FACTORY
LATIN
NORTH
SOUTH
```

## The City or State Key Type File (`bizCityorStateTypeKeys.dat`)

The city or state key type file lists various cities and states that might be used in business names. It also classifies each entry as a city (CT) or state (ST) and indicates the country in which the city or state is located. This enables the Sun Match Engine to recognize and process these values correctly. You can add entries to the city or state key type file using the following syntax.

```
city-or-state type country
```

[Table 25](#) describes the columns in the `bizCityorStateTypeKeys.dat` file.

**TABLE 25** City or State Key Type File

Column	Description
city-or-state	The name of a city or state used in business names.
type	An indicator of whether the value is a city or state. “CT” indicates city and “ST” indicates state.
country	The country code of the country in which the city or state is located.

Following is an excerpt from the `bizCityorStateTypeKeys.dat` file.

```
ADELAIDE          CT  AU
ALABAMA           ST  US
ALASKA            ST  US
ALGIERS           CT  DZ
AMSTERDAM         CT  NL
```

ARIZONA	ST	US
ARKANSAS	ST	US
ASUNCION	CT	PY
ATHENS	CT	GR

## The Business Former Name Reference File (bizCompanyFormerNames.dat)

The business former name reference file provides a list of common company names along with names by which the companies were formerly known so the Sun Match Engine can recognize a business when a record processing a record containing a previous business name. You can add entries to the business former name table using the following syntax.

```
former-name current-name
```

[Table 26](#) describes each column in the bizCompanyFormerNames.dat file.

**TABLE 26** Business Former Name Reference File

Column	Description
former-name	One of the company's previous names.
current-name	The company's current name.

Below is an excerpt from the bizCompanyFormerNames.dat file.

HELLENIC BOTTLING	COCA-COLA HBC
INTERNATIONAL PRODUCTS	THE TERLATO WINE
ORGANIC FOOD PRODUCTS	SPECTRUM ORGANIC PRODUCTS
SUTTER HOME WINERY	TRINCHERO FAMILY ESTATES

## The Merged Business Name Category File (bizCompanyMergerNames.dat)

The merged business name category file provides a list of companies whose name changed because of a merger along with the name of the company after the merge. It also classifies the business names into industry sectors and sub-sectors. This enables the Sun Match Engine to recognize the current company name and determine the sector of the business. You can add entries to the business merger name file using the following syntax.

```
former-name/merged-name sector-code
```

[Table 27](#) describes each column in the bizCompanyMergerNames.dat file.

TABLE 27 Business Merger Name Category File

Column	Description
former-name	The name of the company whose name was not kept after the merger.
merged-name	The name of the company whose name was kept after the merger.
sector-code	The industry sector code of the business. Sector codes are listed in the bizIndustryCategoriesCode.dat file.

Below is an excerpt from the bizCompanyMergerNames.dat file.

```
DUKE/FLUOR DANIEL                20005
FAULTLESS STARCH/BON AMI         09004
FIND/SVP                          10013
FIRST WAVE/NEWPARK SHIPBUILDING  27005
GUNDLE/SLT                       19020
HMG/COURTLAND                   23004
J BROWN/LMC                     10014
KORN/FERRY                      10020
LINSKO/PRIVATE LEDGER           14005
```

## The Primary Business Name Reference File (bizCompanyPrimaryNames.dat)

The primary business name reference file provides a list of companies by their primary name. It also classifies the business names into industry sectors and sub-sectors. This enables the Sun Match Engine to determine the correct value of the sector field when parsing the business name. You can add entries to the primary business name file using the following syntax.

```
primary-name sector-code
```

Table 28 describes the columns in the bizCompanyPrimaryNames.dat file.

TABLE 28 Business Primary Name Reference File

Column	Description
primary-name	The primary name of the company.
sector-code	The industry sector code of the business. Sector codes are listed in the bizIndustryCategoriesCode.dat file.

Below is an excerpt from the bizCompanyPrimaryNames.dat file.

```
BROTHER INTERNATIONAL            12006
BRYSTOL -MYERS SQUIBB           11005
BURLINGTON COAT FACTORY        24003
```

BURLINGTON NORTHERN SANTA FE	27005
BV SOLUTIONS	06012
CABLEVISION	26001
CABOT	04006
CADENCE	06010
CAMPBELL	22006
CAPITAL BLUE CROSS	17001

## The Connector Tokens Reference File (bizConnectorTokens.dat)

The connector tokens reference file defines common values (typically conjunctions) that connect words in business names. For example, in the business name “Nursery of Venice”, “of” is a connector token. This helps the Sun Match Engine recognize and process the full name of a business by indicating that the token connects two parts of the full name.

This file contains one column that lists the connector tokens in the business names you process. You can add entries as needed. Below is an excerpt from the `bizConnectorTokens.dat` file.

```
AN
DE
DES
DOS
LA
LAS
LE
OF
THE
```

## The Country Key Type File (bizCountryTypeKeys.dat)

The country key type file lists countries and continents, along with their abbreviations and assigned nationalities. For continents, the abbreviation is “CON” to separate them from countries. This enables the Sun Match Engine to recognize and process these values as countries or continents. You can add entries to the country key type file using the following syntax.

```
country abbreviation nationality
```

[Table 29](#) describes the columns in the `bizCountryTypeKeys.dat` file.

**TABLE 29** Country Key Type Files

Column	Description
country	The name of a country or continent.
abbreviation	The common abbreviation for the specified country. The abbreviation for a continent is always “CON”.

TABLE 29 Country Key Type Files (Continued)

Column	Description
nationality	The nationality assigned to a person or business originating in the specified country.

Following is an excerpt from the `bizCountryTypeKeys.dat` file.

```

AMERICA                CON  AMERICAN
AFRICA                 CON  AFRICAN
EUROPE                 CON  EUROPEAN
ASIA                   CON  ASIAN
AFGHANISTAN           AF   AFGHAN
ALBANIA                AL   ALBANIAN
ALGERIA                DZ   ALGERIAN

```

## The Industry Sector Reference File (`bizIndustryCategoryCode.dat`)

The industry sector reference file lists and groups various industry sectors and sub-sectors, and includes an identification code for each type so the Sun Match Engine can identify and process the industry sectors for different businesses. You can add entries to the industry sector reference file using the following syntax.

```
sector-code industry-sector
```

Table 30 describes each column in the `bizIndustryCategoryCode.dat` file.

TABLE 30 Industry Sector Reference File

Column	Description
sector-code	The identification code of the specified sector. The first two numbers of each code identify the general industry sector; the last three number identify a sub-sector.
industry-sector	A description of the industry category. This is written in the format “sector - sub-sector”, where sector is a general category of industry types, and sub-sector is a specific industry within that category.

Following is an excerpt from the `bizIndustryCategoryCode.dat` file.

```

02006    Automotive & Transport Equipment - Recreational Vehicles
02007    Automotive & Transport Equipment - Shipbuilding & Related Services
02008    Automotive & Transport Equipment - Trucks, Buses & Other Vehicles
03001    Banking - Banking
04001    Chemicals - Agricultural Chemicals
04002    Chemicals - Basic & Intermediate Chemicals & Petrochemicals
04003    Chemicals - Diversified Chemicals
04004    Chemicals - Paints, Coatings & Other Finishing Products

```

04005	Chemicals - Plastics & Fibers
04006	Chemicals - Specialty Chemicals
05001	Computer Hardware - Computer Peripherals
05002	Computer Hardware - Data Storage Devices
05003	Computer Hardware - Diversified Computer Products

## The Industry Key Type File (bizIndustryTypeKeys.dat)

The industry key type file is used to standardize the value of the Industry field into common industries to which businesses belong so the Sun Match Engine can recognize and process the industry types for different businesses. You can add entries to the industry key type file using the following syntax.

```
industry-type standardized-form sectors
```

Table 31 describes each column in the bizIndustryTypeKeys.dat file.

TABLE 31 Industry Key Type File

Column	Description
industry-type	The original value of the industry type in the input record.
standardized-form	The normalized version of the industry type. If this column contains a name instead of a zero, that name must also be listed in a different entry as an industry type with a standardized form of "0".
sectors	The industry categories of the specified industry type. These values correspond to the sector codes listed in the industry sector file (bizIndustryCategoryCode.dat). You can list as many categories as apply for each type, but they must be entered with a space between each and no line breaks, and they must correspond to an entry in the industry sector file.

Below is an excerpt from the bizIndustryTypeKeys.dat file.

```
TECH          TECHNOLOGY          05001-05007
TECHNOLOGIES TECHNOLOGY          05001-05007
TECHNOLOGY    0                    05001-05007
TECHSYSTEMS   0                    05001-05007
TELE PHONE    TELEPHONE            16005
TELE PHONES   TELEPHONES           16005
TELEVISION    TV                   11013 21014
TELECOM       0                    16005 26006 26009 26010
TELECOMM      TELECOMMUNICATION   16005 26006 26008
TELECOMMUNICATION 0                    16005 26006 26008
```

## The Organization Key Type File (bizOrganizationTypeKeys.dat)

The organization key type file is used to standardize the value of the Organization field into common organizations to which businesses belong. This helps the Sun Match Engine recognize and process the organization types for different businesses. You can add entries to the organization key type file using the following syntax.

```
original-type standardized-form
```

[Table 32](#) describes each column in the bizOrganizationTypeKeys.dat file.

**TABLE 32** Organization Key Type File

Column	Description
original-type	The original value of the organization field in an input record.
standardized-form	The normalized version of an organization type. A zero (0) in this field indicates that the value in the first column is already in its standardized form. If this column contains a name instead of a zero, that name must also be listed in a different entry as an original type with a standardized form of “0”.

Below is an excerpt from the bizOrganizationTypeKeys.dat file.

```
INC                INCORPORATED
INCORPORATED      0
KG                0
KK                0
LIMITED           0
LIMITED PARTNERSHIP 0
LLC               0
LLP               0
LP                LIMITED PARTNERSHIP
LTD               LIMITED
```

## The Business Patterns File (bizPatterns.dat)

The business patterns file defines multiple formats expected from the business name input fields along with the standardized output of each format. The patterns and output appear in two-row pairs in this file, as shown below.

```
4 PNT AST SEP-GLC ORT
PNT AST DEL ORT
```

The first line describes the input pattern and the second describes the output pattern using tokens to denote each component. The supported tokens are described in [“Business Name Tokens” on page 80](#). A number at the beginning of the first line indicates the number of components in the given business name format. You can modify this file using the following syntax.

length input-pattern  
output-pattern

[Table 33](#) lists and describes the syntax components.

**TABLE 33** Business Patterns File Components

Component	Description
length	The number of business name components in the input field.
input-pattern	Tokens that represent a possible input pattern from the unparsed business name fields. Each token represents one component. For more information about address tokens, see <a href="#">“Business Name Tokens” on page 80</a> .
output-pattern	Tokens that represent the output pattern for the specified input pattern. Each token represents one component. For more information about business name tokens, see <a href="#">“Business Name Tokens” on page 80</a> .

Below is an excerpt from the `bizPatterns.dat` file.

```
4 PNT AST SEP-GLC ORT
PNT AST DEL ORT
```

```
4 NFG AJT SEP-GLC ORT
PNT PNT DEL ORT
```

```
4 NF AJT SEP-GLC ORT
PNT PNT DEL ORT
```

```
4 CST IDT NF ORT
PNT PNT PNT ORT
```

```
4 PNT AJT SEP-GLC ORT
PNT PNT DEL ORT
```

## Business Name Tokens

The business patterns file uses tokens to denote different components in a business name, such as the primary name, alias type key, URL, and so on. The file uses one set of tokens for input fields and another set for output fields. The tokens indicate the type key files to use to determine the appropriate values for each output field. You can use only the predefined tokens to represent business name components; the Sun Match Engine does not recognize custom tokens.

[Table 34](#) lists and describes each input token; [Table 35](#) lists and describes each output token.



TABLE 34 Business Name Input Pattern Tokens

Pattern Identifier	Description
CTT	A connector token
PNT	A primary name of a business
PN-PN	A hyphenated primary name of a business
BCT	A common business term
URL	The URL of the business' web site
ALT	A business alias type key (usually an acronym)
CNT	A country name
NAT	A nationality
CST	A city or state type key
IDT	An industry type key
IDT-AJT	Both an industry and an adjective type key
AJT	An adjective type key
AST	An association type key
ORT	An organization type key
SEP	A separator key
NFG	Generic term, not recognized as a specific business name component, with an internal hyphen
NF	Generic term, not recognized as a specific business name component
NFC	A single character, not recognized as a specific business name component
SEP-GLC	A joining comma (a <i>glue type separator</i> )
SEP-GLD	A joining hyphen (a <i>glue type separator</i> )
AND	The text "and"
GLU	A glue type key, such as a forward slash, connecting two parts of a business name component
PN-NF	A business primary name followed by a hyphen and a generic term that is not recognized as a specific business name component
NF-PN	A generic term that is not recognized as a specific business name component, followed by a hyphen and a recognized business primary name

TABLE 34 Business Name Input Pattern Tokens (Continued)

Pattern Identifier	Description
NF-NF	Two generic terms, not recognized as specific business name components and separated by a hyphen

Table 35 lists and describes each output token.

TABLE 35 Business Name Output Pattern Tokens

Pattern Identifier	Description
PNT	The primary name of the business
URL	The URL of the business
ALT	The alias type key of the business (usually an acronym)
IDT	The industry type key of the business
AST	The association type key of the business
ORT	The organization type key of the business
NF	A generic term not recognized as a business name component

## The Special Characters Reference File (bizRemoveSpecChars.dat)

The special characters reference file lists certain characters that should be removed from a business name prior to processing the field, which typically include punctuation marks such as exclamation points, parenthesis, and so on. This enables the Sun Match Engine to recognize the business name.

This file contains one column that lists the characters to be removed from the business names you process. You can add entries as needed. Below is an excerpt from the `bizRemoveSpecChars.dat` file.

```
[
]
{
}
<
>
/
?
```

## Modifying Sun Match Engine Business Name Configuration Files

To customize the Sun Match Engine configuration files for processing business names, you can modify any of the business name standardization files using the text editor provided in NetBeans. Before modifying the match configuration file, review the information provided in [“Sun Match Engine Matching Configuration” on page 12](#) and [“Match Configuration Comparison Functions for Sun Match Engine \(Repository\)” on page 93](#). Make sure a thorough data analysis has been performed to determine the best fields for matching, and the best comparison functions to use for each field.

Updating most standardization files is a straightforward process. Make sure to follow the syntax guidelines provided in [“Sun Match Engine Standardization Configuration for Business Names” on page 69](#). If you add rows to any standardization files, make sure to modify the corresponding parameter in the business constants file (`bizConstants.cfg`). Before making any changes to the patterns file, make sure you understand the tokens used to represent business name field components.

For information about the standardization files you can modify, see [“Sun Match Engine Standardization Configuration for Business Names” on page 69](#).

## Configuring the Matching Service for Business Names (Repository)

To ensure correct processing of business names, you must customize the Matching Service. This includes modifying the Match Field file to support the fields on which you want to match, to standardize the appropriate fields, and to specify the Sun Match Engine as the match and standardization engine (by default, the Sun Match Engine is already specified so this does not need to be changed). Perform the following tasks to configure the Matching Service.

- [“Configuring the Standardization Structure for Business Names \(Repository\)” on page 83](#)
- [“Configuring the Match String for Business Names \(Repository\)” on page 86](#)

When configuring the Matching Service, keep in mind the information presented in [“Configuring the Master Index Matching Service \(Repository\)” on page 27](#).

## Configuring the Standardization Structure for Business Names (Repository)

The standardization structure is configured in the `StandardizationConfig` section of the Match Field file, which is described in detail in [“Match Field Configuration \(Repository\)” in \*Understanding Sun Master Index Configuration Options \(Repository\)\*](#). To configure the required

fields for standardization and phonetic encoding, modify the standardization and phonetic encoding structures. The following sections provide additional guidelines and samples specific to standardizing business names.

---

**Note** – In the default configuration, the rules defined for the business data type assume that all input fields must be parsed as well as normalized. Thus, there is no need to configure fields only for normalization.

---

## Business Name Standardization Structures

For business name fields, the source fields in the standardization structure must include the fields predefined for parsing and normalization. This includes any fields containing business name information, which are parsed into the business name fields listed in “[Business Name Object Structure](#)” on page 68 (excluding the phonetic business name field). The target fields can include any of these parsed fields. Follow the instructions under “Defining Master Index Standardization Rules (Repository)” in *Configuring Sun Master Indexes (Repository)* to define fields for normalization. For the *standardization-type* element, enter **BusinessName** (for more information, see “[Sun Match Engine Match and Standardization Types](#)” on page 24). For a list of field IDs to use in the *standardized-object-field-id* element, see [Table 3](#).

A sample standardization structure for business name data is shown below. This structure parses a business name field into the standard business name fields. Note that there is no domain selector specified, which would normally default to the United States domain; however, since business names are not domain dependent, it is irrelevant here.

```
<free-form-texts-to-standardize>
  <group standardization-type="BusinessName">
    <unstandardized-source-fields>
      <unstandardized-source-field-name>Company.Name
    </unstandardized-source-field-name>
    </unstandardized-source-fields>
    <standardization-targets>
      <target-mapping>
        <standardized-object-field-id>PrimaryName
      </standardized-object-field-id>
      <standardized-target-field-name>Company.Name_Name
      </standardized-target-field-name>
    </target-mapping>
      <target-mapping>
        <standardized-object-field-id>OrgTypekeyword
      </standardized-object-field-id>
      <standardized-target-field-name>Company.Name_OrgType
      </standardized-target-field-name>
    </target-mapping>
  </target-mapping>
</group>
</free-form-texts-to-standardize>
```

```

        <standardized-object-field-id>AssocTypeKeyword
    </standardized-object-field-id>
    <standardized-target-field-name>Company.Name_AssocType
    </standardized-target-field-name>
</target-mapping>
<target-mapping>
    <standardized-object-field-id>IndustrySectorList
    </standardized-object-field-id>
    <standardized-target-field-name>Company.Name_Sector
    </standardized-target-field-name>
</target-mapping>
<target-mapping>
    <standardized-object-field-id>IndustryTypeKeyword
    </standardized-object-field-id>
    <standardized-target-field-name>Company.Name_Industry
    </standardized-target-field-name>
</target-mapping>
<target-mapping>
    <standardized-object-field-id>AliasList
    </standardized-object-field-id>
    <standardized-target-field-name>Company.Name_Alias
    </standardized-target-field-name>
</target-mapping>
<target-mapping>
    <standardized-object-field-id>Url
    </standardized-object-field-id>
    <standardized-target-field-name>Company.Name_URL
    </standardized-target-field-name>
</target-mapping>
</standardization-targets>
</group>
</free-form-texts-to-standardize>

```

## Business Name Phonetic Encoding

When you match on business name fields, the name field should be specified for phonetic conversion (by default, the wizard defines this for you). Follow the instructions under “Defining Phonetic Encoding for the Master Index (Repository)” in *Configuring Sun Master Indexes (Repository)* to define fields for phonetic encoding.

A sample of the *phoneticize-fields* element is shown below. This sample only converts the business name. You can define additional fields for phonetic encoding.

```

<phoneticize-fields>
  <phoneticize-field>
    <unphoneticized-source-field-name>Company.Name_Name
    </unphoneticized-source-field-name>
    <phoneticized-target-field-name>Company.Name_NamePhon

```

```
</phoneticized-target-field-name>  
<encoding-type>NYSIIS</encoding-type>  
</phoneticize-field>  
</phoneticize-fields>
```

## Configuring the Match String for Business Names (Repository)

For matching on business name fields, make sure the match string you specify in the MatchingConfig section of the Match Field file contains all or a subset of the fields that contain the standardized data (the unparsed business names are typically too inconsistent for matching). You can include additional fields for matching if required.

To configure the match string, follow the instructions under “Defining the Master Index Match String (Repository)” in *Configuring Sun Master Indexes (Repository)*. For the Sun Match Engine, each data type has a different match type (specified by the *match-type* element). The PrimaryName, OrgTypeKeyword, AssocTypeKeyword, IndustrySectorList, IndustryTypeKeyword, and Url match types are specific to business name matching. You can specify any of the other match types defined in the match configuration file, as well. For more information, see “[Sun Match Engine Match and Standardization Types](#)” on page 24.

A sample match string for business name matching is shown below. This sample matches on the company name, the organization type, and the sector.

```
<match-system-object>  
  <object-name>Company/object-name</object-name>  
  <match-columns>  
    <match-column>  
      <column-name>Enterprise.SystemSBR.Company.Name_PrimaryName</column-name>  
      <match-type>PrimaryName</match-type>  
    </match-column>  
    <match-column>  
      <column-name>Enterprise.SystemSBR.Company.Name_OrgType</column-name>  
      <match-type>OrgTypeKeyword</match-type>  
    </match-column>  
    <match-column>  
      <column-name>Enterprise.SystemSBR.Company.Name_Sector</column-name>  
      <match-type>IndustryTypeKeyword</match-type>  
    </match-column>  
  </match-columns>  
</match-system-object>
```

# Fine-Tuning Weights and Thresholds for Sun Match Engine (Repository)

Each Sun Match Engine implementation is unique, typically requiring extensive data analysis to determine how to best configure the structure and matching logic of the master index application. The following topics provide an overview of the process of fine-tuning the matching logic in the match configuration file and fine-tuning the match and duplicate thresholds.

- [“Data Analysis Overview” on page 87](#)
- [“Customizing the Match Configuration and Thresholds” on page 87](#)

## Data Analysis Overview

A thorough analysis of the data to be shared with the master index application is a must before beginning any implementation. This analysis not only defines the types of data to include in the object structure, but indicates the relative reliability of each system’s data, helps determine which fields to use for matching, and indicates the relative reliability of each match field.

To begin the analysis, the legacy data that will be converted into the master index database is extracted and analyzed. Once the initial analysis is complete, you can perform an iterative process to fine-tune the matching and duplicate thresholds and to determine the level of potential duplication in the existing data. If you plan to use the Data Profiler and Bulk Matcher tools generated by Sun Match Engine to analyze data, review the information in *Analyzing and Cleansing Data for Sun Master Index* before you extract the legacy data.

---

**Note** – These tools are only available from service-enabled master index applications. The document referenced above describes what you need to do to generate the tools if you are using Sun Match Engine (Repository).

---

## Customizing the Match Configuration and Thresholds

There are three primary steps to customizing how records are matched in a master index application.

- [“Determining the Match Fields” on page 88](#)
- [“Customizing the Match Configuration” on page 88](#)
- [“Determining the Weight Thresholds” on page 90](#)

## Determining the Match Fields

Before extracting data for analysis, review the types of data stored in the messages generated by each system. Use these messages to determine which fields and objects to include in the object structure of the master index application. From this object structure, select the fields to use for matching. When selecting these fields, keep in mind how representative each field is of a specific object. For example, in a master person index, the social security number field, first and last name fields, and birth date are good representations whereas marital status, suffix, and title are not. Certain address information or a home telephone number might also be considered. In a master company index, the match fields might include any of the fields parsed from the complete company name field, as well as a tax ID number or address and telephone information.

## Customizing the Match Configuration

Once you determine the fields to use for matching, determine how the weights will be generated for each field. The primary tasks include determining whether to use probabilities or agreement weight ranges and then choosing the best comparison functions to use for each match field.

### Probabilities or Agreement Weights

The first step in configuring the match configuration is to decide whether to use m-probabilities and u-probabilities or agreement and disagreement weight ranges. Both methods will give you similar results, but agreement and disagreement weight ranges allow you to specify the precise maximum and minimum weights that can be applied to each match field, giving you control over the value of the highest and lowest matching weights that can be assigned to each record.

### Defining Relative Value

For each field used for matching, define either the m-probabilities and u-probabilities or the agreement and disagreement weight ranges in the match configuration file. Review the information provided under [“Sun Match Engine Matching Weight Formulation” on page 9](#) to help determine how to configure these values. Remember that a higher m-probability or agreement weight gives the field a higher weight when field values agree.

### Determining the Weight Range

In order to find the initial values to set for the match and duplicate thresholds, you must determine the total range of matching weights that can be assigned to a record. This weight is the sum of all weights assigned to each match field. Running the Bulk Matcher in match analysis mode can help you determine the match and duplicate thresholds. For more information about this tool, see “Performing a Match Analysis” in *Loading the Initial Data Set for a Sun Master Index*.

The way you determine weight ranges varies depending on whether you are using m and u-probabilities or agreement and disagreement weights.



## Weight Ranges Using Agreement Weights

For agreement and disagreement weight ranges, determining the match weight ranges is very straightforward. Simply total the maximum agreement weights for each field to determine the maximum match weight. Then total the minimum disagreement weights for each match field to determine the minimum match weight. [Table 36](#) provides a sample agreement/disagreement configuration for matching on person data. As you can see, the range of match weights generated for a master index application with this configuration is from -36 to +38.

**TABLE 36** Sample Agreement and Disagreement Weight Ranges

Field Name	Maximum Agreement Weight	Minimum Disagreement Weight
First Name	8	-8
Last Name	8	-8
Date of Birth	7	-5
Gender	5	-5
SSN	10	-10
<b>Maximum Match Weight</b>	38	
<b>Minimum Match Weight</b>		-36

## Weight Ranges Using Probabilities

Determining the match weight ranges when using m-probabilities and u-probabilities is a little more complicated than using agreement and disagreement weights. To determine the maximum weight that will be generated for each field, use the following formula:

$$\text{LOG2}(m\_prob/u\_prob)$$

To determine the minimum match weight that will be generated for each field, use the following formula:

$$\text{LOG2}((1-m\_prob)/(1-u\_prob))$$

[Table 37](#) below illustrates a sample of m-probabilities and u-probabilities, including the corresponding agreement and disagreement weights that are generated with each combination of probabilities. As you can see, the range of match weights generated for a master index application with this configuration is from -35.93 to +38

TABLE 37 Sample m-probabilities and u-probabilities

Field Name	m-probability	u-probability	Max Agreement Weight	Min Disagreement Weight
First Name	.996	.004	7.96	-7.96
Last Name	.996	.004	7.96	-7.96
Date of Birth	.97	.007	7.11	-5.04
Gender	.97	.03	5.01	-5.01
SSN	.999	.001	9.96	-9.96
<b>Maximum Match Weight</b>			38	
<b>Minimum Match Weight</b>				-35.93

## Comparison Functions

The match configuration file defines several match types for different types of fields. You can either modify existing rows in this file or create new rows that define custom matching logic. To determine which comparison functions to use, review the information provided in [“Match Configuration Comparison Functions for Sun Match Engine \(Repository\)”](#) on page 93. Choose the comparison functions that best suit how you want the match fields to be processed.

## Determining the Weight Thresholds

Weight thresholds tell the master index application how to process incoming records based on the matching probability weights generated by the Sun Match Engine. Two parameters in the Threshold file provide the master index application with the information needed to determine if records should be flagged as potential duplicates, if records should be automatically matched, or if a record is not a potential match to any existing records.

- **Match Threshold** - Specifies the weight at which two profiles are assumed to represent the same person and are automatically matched (this depends on the setting for the OneExactMatch parameter).
- **Duplicate Threshold** - Specifies the minimum weight at which two profiles are considered potential duplicates of one another. The matching threshold indicates the maximum weight for potential duplicates.

[Figure 1](#) illustrates the match and duplicate thresholds in comparison to total composite match weights.

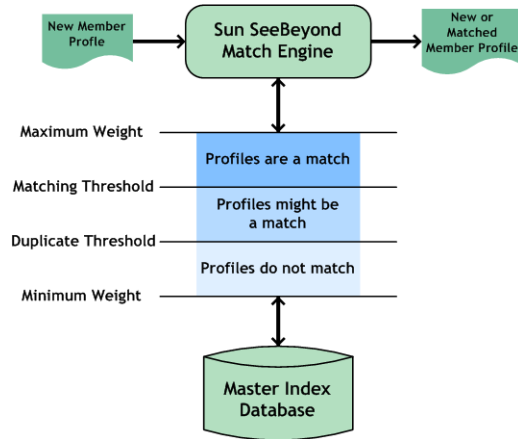


FIGURE 1 Weight Thresholds

## Specifying the Weight Thresholds

There are many techniques for determining the initial settings for the match and duplicate thresholds. This section discusses two methods.

The first method, the weight distribution method, is based on the calculation of the error rates of false matches and false non-matches from analyzing the distribution spectrum of all the weighted pairs. This is the standard method, and is illustrated in [Figure 2](#). The second method, the percentage method relies on measuring the total maximum and minimum weights of all the matched fields and then specifying a certain percentage of these values as the initial thresholds.

The weight distribution method is more thorough and powerful but requires analyzing a large amount of data (match weights) to be statistically reliable. It does not apply well in cases where one candidate record is matched against very few reference records. The percentage method, though simple, is very reliable and precise when dealing with such situations. For both methods, defining the match threshold and the duplicate threshold is an iterative process.

### Weight Distribution Method

Each record pair in the master index application can be classified into three categories: matches, non-matches, and potential matches. In general, the distribution of records is similar to the graph shown in [Figure 2](#). Your goal is to make sure that very few records fall into the False Matches region (if any), and that as few as possible fall into the False Non-matches region. You can see how modifying the thresholds changes this distribution. Balance this against the number of records falling within the Manual Review section, as these will each need to be reviewed, researched, and resolved individually.

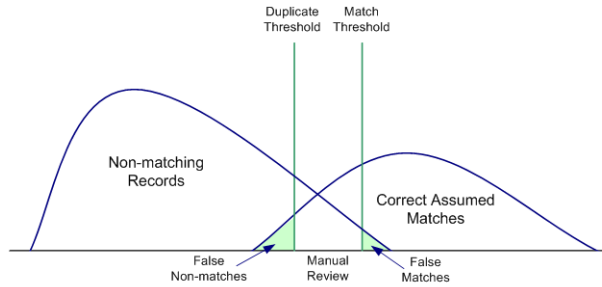


FIGURE 2 Weight Distribution Chart

### Percentage Method

Using this method, you set the initial thresholds as a percentage of the maximum and minimum weights. Using the information provided under [“Weight Ranges Using Agreement Weights” on page 89](#) or [“Weight Ranges Using Probabilities” on page 89](#), determine the maximum and minimum values that can be generated for composite match weights. For the initial run, the match threshold is set intentionally high to catch only the most probable matches. The duplicate threshold is set intentionally low to catch a large set of possible matches.

Set the match threshold at 70% of the maximum composite weight starting from zero as the neutral value. Using the weight range samples in [Table 37](#), this would be 70% of 38, or 26.6. Set the duplicate threshold near the neutral value (that is, the value in the center of the maximum and minimum weight range). The value could be set between 10% of the maximum weight and 10% of the minimum weight. Using the samples above, this would be between 3.8 (10% of 38) and -3.6 (10% of -36).

### Fine-tuning the Thresholds

Achieving the correct thresholds for your implementation is an iterative process. First, using the initial thresholds, process the data extracts into the master index database. Then analyze the resulting assumed match and potential duplicates, paying close attention to the assumed match records with matching weights close to the match threshold, to potential duplicate records close to either threshold, and to non-matches near the duplicate threshold.

If you find that most or all of the assumed matches at the low end of the match range are not actually duplicate records, raise the match threshold accordingly. If, on the other hand, you find several potential duplicates at the high end of the duplicate range that are actual matches, decrease the match threshold accordingly. If you find that most or all of the potential duplicate records in the low end of the duplicate range should not be considered duplicate matches, consider raising the duplicate threshold. Conversely, if you find several non-matches with weight near the duplicate threshold that should be considered potential duplicates, lower the duplicate threshold.

Repeat the process of loading and analyzing data and adjusting the thresholds until you are satisfied with the results.

## Match Configuration Comparison Functions for Sun Match Engine (Repository)

Match field comparison functions, or *comparators*, compare the values of a field in two records to determine whether the fields match or how closely they match. The fields are then assigned a matching weight based on the results of the comparison function. You can use several different types of comparison functions in the match configuration file in order to customize how the Sun Match Engine matches records. The following topics provide information about the predefined comparison functions and their parameters and options.

- [“Sun Match Engine Comparison Functions” on page 93](#)
- [“Sun Match Engine Comparison Function Options” on page 103](#)

### Sun Match Engine Comparison Functions

There are six primary types of comparison functions used by the Sun Match Engine. The following types of comparison functions are available.

- [“Bigram Comparators” on page 93](#)
- [“Uncertainty String Comparators” on page 94](#)
- [“Exact char-by-char Comparator \(c\)” on page 97](#)
- [“Numeric Comparators” on page 97](#)
- [“Date Comparators” on page 99](#)
- [“Prorated Comparator \(p\)” on page 101](#)

Certain comparison function types are very specific to the type of data being matched, such as the numeric functions and the date functions. Others, such as the Bigram and uncertainty functions, are more general and can be applied to various data fields.

Be sure to review [Table 1](#) for information about how the parameters in the match configuration file affect the outcome of the comparator functions. For example, these parameters define how null fields are handled and what the actual agreement and disagreement weights will be.

---

**Note** – The names of the comparators are configurable. The default names are used here.

---

### Bigram Comparators

The Sun Match Engine provides two different comparison functions based on the Bigram algorithm, the standard bigram (b1) and the transposition bigram (b2). A Bigram algorithm compares two strings using all combinations of two consecutive characters within each string.

For example, the word “bigram” contains the following bigrams: “bi”, “ig”, “gr”, “ra”, and “am”. The Bigram comparison function returns a value between 0 and 1, which accounts for the total number of bigrams that are in common between the strings divided by the average number of bigrams in the strings. Bigrams handle minor typographical errors well.

### **Bigram String Comparator (b1)**

This is a standard Bigram comparison function, processing match fields as described above. This comparison function takes no parameters.

### **Advanced Bigram String Comparator (b2)**

This comparison function is based on the standard Bigram comparison function, but handles transpositions of characters within a string. This comparison function takes no parameters.

## **Uncertainty String Comparators**

The Sun Match Engine provides the following uncertainty comparison functions for comparing string fields. Most uncertainty comparison functions are generic, but three comparison functions are designed for specific types of information (first name, last name, and house number).

- “Generic String Comparator (u)” on page 94
- “Advanced Generic String Comparator (ua)” on page 95
- “Simplified String Comparator (us)” on page 95
  - “Simplified String Comparator - FirstName (uf)” on page 95
  - “Simplified String Comparator - LastName (ul)” on page 96
  - “Simplified String Comparator - House Numbers (un)” on page 96
- “Language-specific String Comparator (usu)” on page 96

### **Generic String Comparator (u)**

This is the standard uncertainty comparison function, which processes string fields as described above. As more differences are found between two fields, the agreement weight decreases nonlinearly. Thus, the agreement weight can remain high for several differences, but will drop sharply at a certain point. This comparison function takes no parameters.

The uncertainty comparison function is based on the Jaro algorithm with McLaughlin adjustments for similarities. The Jaro algorithm is a string comparison function that accounts for insertions, deletions, and transpositions by performing the following steps.

1. Compute the lengths of both strings to be matched.
2. Determine the number of common characters between the two strings. In order for characters to be considered common, they must be within one-half the length of the shorter string.

3. Determine the number of transpositions. A transposition means a character from the first string is out of order with the corresponding common character from the second string.

### Advanced Generic String Comparator (ua)

This comparison function is based on the standard uncertainty comparison function, **u**, with variants of Winkler/Lynch and McLaughlin. It has additional features to handle specific differences between fields, such as key punch and visual memory errors. Each feature makes use of the information made available from previous features. This comparison function takes no parameters. The following features are included in the advanced uncertainty function.

- The function determines each character in exact agreement and then assigns a value of 1.0 to each agreeing character. It then determines each disagreeing but similar character and assigns a value of 0.3 to each. Similar characters might occur because of scanning errors (for example, “1” the number versus “l” the letter) or keypunch errors (for example, “S” versus “D”).
- The function gives increased value to agreement on the beginning characters of a string. The algorithm adjusts the weighting value up by a fixed amount if the first four characters in each string agree; it adjusts the weighting value up by smaller value if only the first three, two, or one characters agree.
- The function adjusts the string comparison value if the strings are longer than six characters and more than half of the characters after the fourth character agree.

### Simplified String Comparator (us)

This comparison function is a custom version of a generic string comparison function. It is similar to the basic uncertainty comparison function, **u**, but processes data in a more simple and efficient manner, improving processing speed. The agreement weights generated by this comparison function decrease in a more uniform manner for each difference found between two fields.

Like the basic uncertainty function, the simplex function takes into account such uncertainty factors as string length, transpositions, key punch errors, and visual memory errors. Unlike the uncertainty comparison function (“u”), this function handles diacritical marks. This comparison function takes no parameters.

### Simplified String Comparator - FirstName (uf)

This comparison function is designed specifically for matching on first name fields, and is based on the simplex uncertainty comparison function, **us**. This comparison function analyzes the string and then adjusts the weight based on statistical data. This comparison function takes no parameters.

## Simplified String Comparator - LastName (ul)

This comparison function is designed specifically for matching on last name fields, and is based on the simplex uncertainty comparison function, **us**. This comparison function analyzes the string and then adjusts the weight based on statistical data. This comparison function takes no parameters.

## Simplified String Comparator - House Numbers (un)

This comparison function is designed specifically for matching on house numbers, and is based on the simplex uncertainty comparison function, **u**. This comparison function analyzes the string and then adjusts the weight based on statistical data. This comparison function takes no parameters.

## Language-specific String Comparator (usu)

This comparison function is a custom version of a generic string comparison function. It is similar to the simplex uncertainty comparison function, **us**, but is based in Unicode to enable multilingual support. This locale-oriented comparator recognizes the nuances of each language and supports the complexities and subtleties of each. For example, when configured to use the German language set, the function recognizes “ß” and “ss” as equivalent. Like the simplex uncertainty function, the Unicode function takes into account such uncertainty factors as string length, transpositions, key punch errors, and visual memory errors. This comparison function takes the parameter described in [Table 38](#).

TABLE 38 usu Comparison Function Parameter

Parameter	Description
language	<p>An indicator of the language being used for the information stored in the database. Enter one of the following codes to indicate the language in use.</p> <p><b>da</b> - Danish</p> <p><b>sv</b> - Swedish</p> <p><b>nb</b> - Norwegian Bokmål</p> <p><b>nn</b> - Norwegian Nynorsk</p> <p><b>nl</b> - Dutch</p> <p><b>es</b> - Spanish</p> <p><b>fr</b> - French</p> <p><b>en</b> - English</p> <p><b>it</b> - Italian</p> <p><b>de</b> - German</p>



## Exact char-by-char Comparator (c)

The Sun Match Engine provides one exact-match comparison function, “c”. With this comparison function, two fields must match on each character in order to be considered a match. This comparison function takes no parameters.

## Numeric Comparators

The Sun Match Engine provides several comparison functions for matching on numeric fields.

- “Generic Number Comparator (n)” on page 98
- “Integer Comparator (nI)” on page 98
- “Real Number Comparator (nR)” on page 98
- “Alphanumeric Comparator (nS)” on page 98

All but the nS comparison function can perform numeric string comparisons or relative distance calculations. When set for a string comparison, the functions compare numeric strings based on the advanced uncertainty comparator. When set for relative distance calculations, the matching weight between two numbers decreases as the numbers become further apart, until the relative distance plus one is reached. At this point, the numbers are considered non-matches. For example, if the relative distance is “10” and the base number for comparison is “2”, a field value of 8 receives a lower matching weight than a field value of 4; but a field value of 13 is considered a complete non-match (since the distance between 2 and 13 is 11).

Figure 3 illustrates how the weight is decreased as the difference between the two compared fields reaches the relative distance. In this diagram, the relative distance is 10 and the light blue line represents the agreement weight. When the difference between two fields reaches 11 (relative distance plus one), the fields are considered a non-match and are given the full disagreement weight.

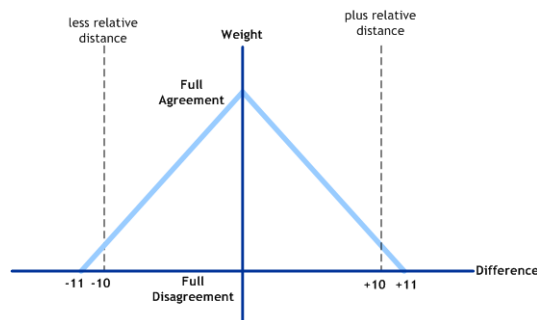


FIGURE 3 Numeric Relative Distance Comparison

## Generic Number Comparator (n)

This is a basic numeric comparison function, processing numeric fields as described above. It accepts the parameters listed in [Table 39](#).

TABLE 39 n, nI, and nR Comparison Function Parameters

Parameter	Description
distance-or-string	Specifies whether a relative distance calculation or a direct string comparison is used. Specify “y” to use a relative distance calculation; specify “n” to use a string comparison.
relative-distance	The greatest difference between two integers at which the values could still be considered a possible match. When the difference between two numbers is greater than the relative distance, the numbers are considered a non-match (the weight becomes zero when the actual difference is the relative distance plus one).

## Integer Comparator (nI)

This numeric comparison function matches specifically on integers and accepts the parameters listed in [Table 39](#).

## Real Number Comparator (nR)

This numeric comparison function matches specifically on real numbers and accepts the parameters listed in [Table 39](#).

## Alphanumeric Comparator (nS)

This numeric comparison function is designed specifically for matching on numeric strings and is very useful for matching social security numbers or other unique identifiers. This is the only numeric comparator that can compare alphanumeric values rather than just numeric values. It accepts the parameters listed in [Table 40](#).

TABLE 40 nS Comparison Function Parameters

Parameter	Description
fixed-length	An optional parameter that takes the length of the field value into account. If a fixed length is specified, the match engine considers any field of a different length to be a non-match. Specify any integer smaller than the value specified for the size specified for the field (for more information, see <a href="#">“Matching Rules” on page 13</a> ).
character-type	An indicator of whether the field must be all numeric. Specify “nu” for numeric only, or specify “an” to allow alphanumeric characters. The match engine considers any fields containing characters that are not allowed to be a non-match.

TABLE 40 nS Comparison Function Parameters (Continued)

Parameter	Description
invalid-characters	A list of invalid characters for the field. If you specify a character, the match engine considers fields that consist of only that character to be a non-match. For example, if you specify “0”, then an SSN field cannot contain all zeros. Specify as many alphanumeric characters as needed, separated by a space.

## Date Comparators

The Sun Match Engine provides various date comparison functions. When comparing dates, the match engine compares each date component (for example, it compares the year in the first date against the year in the second date, the month against the month, and the day against the day). This allows for multiple transpositions in each date field. The date comparators use the Java date format (`java.sql.Date`), allowing the comparator to use the Gregorian calendar and to take into account the time zone where the date field originated.

The following comparison functions are available for matching on date fields.

- [“Date Comparator - Year only \(dY\)” on page 100](#)
- [“Date Comparator - Month-Year \(dM\)” on page 100](#)
- [“Date Comparator - Day-Month-Year \(dD\)” on page 100](#)
- [“Date Comparator - Hour-Day-Month-Year \(dH\)” on page 100](#)
- [“Date Comparator - Min-Hour-Day- Month-Year \(dm\)” on page 101](#)
- [“Date Comparator - Sec-Min-Hour-Day- Month-Year \(ds\)” on page 101](#)

As with the numeric comparison functions, the date comparison functions can use either a direct string comparison or a relative distance calculation. When using a relative distance calculation, the matching weight between two dates decreases as the dates become further apart, until the relative distance is reached. When the difference becomes the relative distance plus one, the dates are considered non-matches. You can specify different relative distances for before and after the given date. Any dates falling outside of the specified time period receive a complete disagreement weight. The relative distances are specified in the smallest unit of time being matched.

[Figure 4](#) illustrates how the weight is decreased as the difference between the two compared fields reaches either the before or after relative distance. In this diagram, the before relative distance is **11**, the after relative distance is **5**, and the light blue line represents the agreement weight. When the base date is later than the compared date and the difference between the dates reaches **11** (distance before plus one), the fields are considered a non-match and are given the full disagreement weight. When the base date is earlier than the compared date and the difference between the dates reaches **6** (distance after plus 1), the fields are considered a non-match.

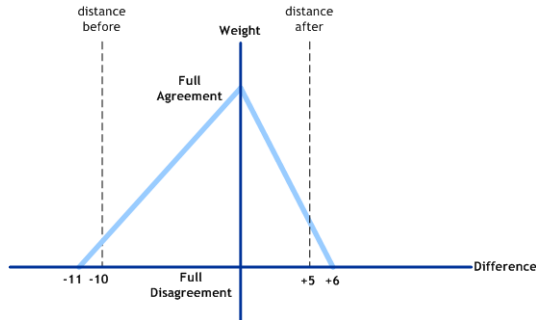


FIGURE 4 Date Relative Distance Comparison

The date comparison functions take the parameters listed in [Table 41](#).

TABLE 41 Date Comparison Function Parameters

Parameter	Description
distance-or-string	Specifies whether a relative distance calculation or a direct string comparison is used. Specify “y” to use a relative distance calculation; specify “n” to use a string comparison.
distance-before	The number of units prior to the reference date/time for which two date fields can still be considered a match.
distance-after	The number of units following the reference date/time for which two date fields can still be considered a match.

### Date Comparator - Year only (dY)

This date comparison function takes only the 4-character year into account for matching. If relative distance calculation is specified, the relative distance is specified in years.

### Date Comparator - Month-Year (dM)

This date comparison function takes the month and year into account for matching. If relative distance calculation is specified, the relative distance is specified in months.

### Date Comparator - Day-Month-Year (dD)

This date comparison function takes the day, month, and year into account for matching. If relative distance calculation is specified, the relative distance is specified in days.

### Date Comparator - Hour-Day-Month-Year (dH)

This date comparison function takes the hour, day, month, and year into account for matching. If relative distance calculation is specified, the relative distance is specified in hours.

### **Date Comparator - Min-Hour-Day- Month-Year (dm)**

This date comparison function takes the minute, hour, day, month, and year into account for matching. If relative distance calculation is specified, the relative distance is specified in minutes.

### **Date Comparator - Sec-Min-Hour-Day- Month-Year (ds)**

This date comparison function takes the second, minute, hour, day, month, and year into account for matching. If relative distance calculation is specified, the relative distance is specified in seconds.

### **Prorated Comparator (p)**

The prorated comparison function uses a relative distance calculation and allows you to specify how quickly the agreement weight between two fields decreases. Matching weights are assigned with a linear adjustment according to the parameters you specify. You specify an initial agreement range. If the difference between two fields falls within that range, the fields are considered a complete match. You also specify a disagreement range ending with the relative distance. If the difference between two fields falls within that range, the fields are considered a non-match. When the difference between the fields falls between those two ranges, they are considered to be partial matches and the agreement weight is adjusted linearly. Any difference greater than the relative distance is always considered a non-match.

[Figure 5](#) illustrates how weighting is adjusted per the parameters you define. In these diagrams, the green line indicates full agreement, the light blue line indicates prorated agreement, and the red line indicates full disagreement. The diagrams illustrate how increasing the disagreement weight causes the prorated agreement weight to decrease more sharply.

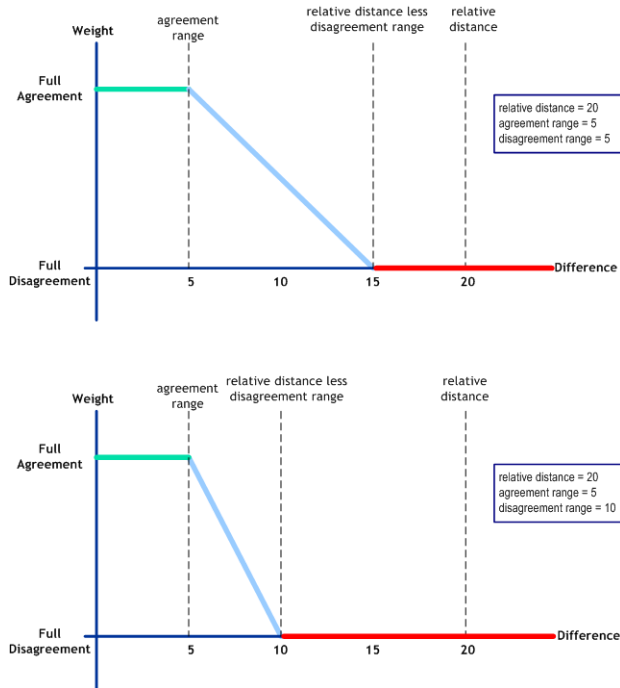


FIGURE 5 Prorated Linear Adjustment Comparison

The prorated comparison function takes the parameters listed in Table 42.

TABLE 42 Prorated Comparison Function Parameters

Parameter	Description
relative-distance	The greatest difference between two numbers at which they can still be considered a match or partial match.
agreement-range	The greatest difference between two numbers at which they are considered a full match. This number must be less than the relative distance.
disagreement-range	This number indicates the minimum difference at which two numbers are considered a non-match and shortens or lengthens the weighting scale. To find this difference, the match engine subtracts this value from the relative distance. If the fields differ by that amount or greater, they are considered to be a non-match.  The weighting scale decreases in size as the value of the full-disagreement parameter increases (see diagram).

## Sun Match Engine Comparison Function Options

The options listed below can be used in conjunction with the above string comparison functions to give them more functionality. For example, you can use an "ufl" string comparison function that refers to the first name comparison function with the possibility to switch fields if the first one does not match.

- **I** - This is a major inversion option that allows for field transpositions. If two compared fields do not match, this option lets the match engine know to compare the original field in the first record with the next field in the second record. If those fields agree, the match engine assigns the full agreement weight and switches the fields.
- **i** - This is a minor inversion option similar to the major inversion (I) described above. This option only assigns one-half of the full agreement weight if the transposed fields match.
- **x** - If two or more fields with this option match, their weight is doubled; but if any of the fields with this option disagree, the weight is not doubled.
- **k** - This option can be used with the "x" option to give more importance to one field. Specifying this option on a field tells the match engine to double the match weight for a sub-group of fields with the "x" option by doubling the weight as soon as it comes to the field with the "k" option.

