



Master Index Match Engine Reference



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

Part No: 821-0919-10
December 2009

Copyright 2009 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 SunTM 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

Master Index Match Engine Reference	7
About the Master Index Match Engine	7
Related Topics	8
Master Index Match Engine Overview	8
Data Matching Concepts	8
Deterministic and Probabilistic Data Matching	9
Weighting Thresholds	9
Probabilities and Direct Weights	9
How the Master Index Match Engine Works	10
Master Index Match Engine Structure	11
Master Index Match Engine Configuration Files	11
Master Index Match Engine Matching Weight Formulation	11
Master Index Match Engine Data Types	12
The Master Index Match Engine and the Master Index Standardization Engine	12
Sun Master Index Standardization and Matching Process	13
Master Index Match Engine Matching Configuration	14
The Master Index Match Engine Match Configuration File	14
Master Index Match Engine Match Configuration File Format	15
Master Index Match Engine Matching Comparison Functions At a Glance	18
Master Index Match Engine Comparator Definition List	21
Master Index Match Engine Comparison Functions	24
Bigram Comparators	25
Bigram Comparator (b1)	25
Advanced Bigram Comparator (b2)	25
Uncertainty String Comparators	25
Advanced Jaro String Comparator (u)	26
Winkler-Jaro String Comparator (ua)	26
Condensed String Comparator (us)	27

Advanced Jaro Adjusted for First Names (uf)	27
Advanced Jaro Adjusted for Last Names (ul)	27
Advanced Jaro Adjusted for House Numbers (un)	27
Advanced Jaro AlphaNumeric Comparator (ujjs)	27
Unicode String Comparator (usu)	28
Unicode AlphaNumeric Comparator (usus)	29
Exact Character-to-Character Comparator (c)	30
Numeric Comparators	31
Integer Comparator (nI)	31
Real Number Comparator (nR)	32
Condensed AlphaNumeric SSN Comparator (nS)	32
Date Comparators	33
Date Comparator With Years as Units (dY)	34
Date Comparator With Months as Units (dM)	35
Date Comparator With Days as Units (dD)	35
Date Comparator With Hours as Units (dH)	35
Date Comparator With Minutes as Units (dm)	35
Date Comparator With Seconds as Units (ds)	35
Prorated Comparator (p)	35
Creating Custom Comparators for the Master Index Match Engine	37
Custom Comparator Overview	37
About the Comparator Package	38
Defining Custom Comparators	39
Before You Begin	39
Step 1: Create the Custom Comparator Java Class	40
Step 2: Register the Comparator in the Comparators List	43
Step 3: Define Parameter Validations (Optional)	45
Step 4: Define Data Source Handling (Optional)	46
DataSourcesProperties Class	47
Step 5: Define Curve Adjustment or Linear Fitting (Optional)	50
Step 6: Compile and Package the Comparator	51
Step 7: Import the Comparator Package Into Sun Master Index	52
Step 8: Configure the Comparator in the Match Configuration File	52
Master Index Match Engine Configuration for Common Data Types	53
The Master Index Match String	53
Master Index Match Engine Match String Fields	54

Person Data Match String Fields	54
Address Data Match String Fields	54
Business Name Match String Fields	55
Master Index Match Engine Match Types	55
Configuring the Match String for a Master Index Application	57
Configuring the Match String for Person Data	58
Configuring the Match String for Address Data	59
Configuring the Match String for Business Names	60
Fine-Tuning Weights and Thresholds for Sun Master Index	61
Data Analysis Overview	61
Customizing the Match Configuration and Thresholds	61
Determining the Match Fields	61
Customizing the Match Configuration	62
Determining the Weight Thresholds	64

Master Index Match Engine Reference

The topics listed here provide conceptual information about the Master Index Match Engine and how it matches data in a master index application. The Master Index Match Engine is designed to work closely with the Master Index Standardization Engine. For more information about the standardization engine, see [Master Index Standardization Engine Reference](#).

- [Master Index Match Engine Overview](#)
- [Master Index Match Engine Matching Configuration](#)
- [Master Index Match Engine Comparison Functions and Options](#)
- [Creating Custom Comparators for the Master Index Match Engine](#)
- [Master Index Match Engine Configuration for Common Data Types](#)
- [Fine-Tuning Weights and Thresholds for Sun Master Index](#)

About the Master Index Match Engine

The Master Index Match Engine provides record matching capabilities for external applications, such as master index applications. It works best along with the Master Index Standardization Engine, which provides the preprocessing of data that is required for accurate matching, such as data parsing, data standardization, and phonetic encoding. 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 the fields on which matching is performed and how the matching logic is configured. The composite weight is usually the sum of weights generated for all match fields in the records (but could also be a function of the match field weights). The composite weight indicates how closely two records match.

Related Topics

Several topics provide information and instructions for implementing and using a master index application. For a complete list of topics related to working with the service-enabled version of Sun Master Index, see [“Related Topics”](#) in *Sun Master Index User’s Guide*.

Master Index Match Engine Overview

The Master Index Match Engine is the standard match engine designed to work with the master index applications created by Sun Master Index. The match engine can also be called from other applications. It is highly configurable in the Sun Master Index environment and can be used to match on various types of data. The match engine is highly configurable in the Sun Master Index environment and can be used to match various types of data. The Master Index Match Engine works in conjunction with the Master Index Standardization Engine to improve the quality of your data.

The following topics provide information about matching concepts, the match process, and how the Master Index Match Engine matches data.

- [“Data Matching Concepts”](#) on page 8
- [“How the Master Index Match Engine Works”](#) on page 10
- [“Sun Master Index Standardization and Matching Process”](#) on page 13

Data Matching Concepts

Data matching compares data stored in disparate systems in and across organizations, helping you reduce data duplication and improve data accuracy. Matching involves comparing specific fields in two standardized records and returning a weight that indicates the likelihood of a match between the two records. A higher weight between two records indicates a greater likelihood of a match. Data matching is based on proven algorithms that are designed to compare different types of data, such as strings, dates, integers, and so on. Matching is a key step in managing data quality, and the algorithms are typically quite complex. Some algorithms are configured to compare more specialized types of data, including first and last names, social security numbers, and dates of various formats.

The following topics provide additional information about standard data matching concepts:

- [“Deterministic and Probabilistic Data Matching”](#) on page 9
- [“Weighting Thresholds”](#) on page 9
- [“Probabilities and Direct Weights”](#) on page 9

Deterministic and Probabilistic Data Matching

Data matching can be either deterministic or probabilistic. In deterministic matching, either unique identifiers for each record are compared to determine a match or an exact comparison is used between fields. Unique identifiers can include national IDs, system IDs, and so on. This can include system IDs, national IDs, and so on. Deterministic matching is generally not completely reliable since in some cases no single field can provide a reliable match between two records. This is where probabilistic, or *fuzzy*, matching comes in. In probabilistic matching, several field values are compared between two records and each field is assigned a weight that indicates how closely the two field values match. The sum of the individual fields weights indicates the likelihood of a match between two records.

Weighting Thresholds

In a data management system, you can set duplicate and match threshold weights. The duplicate threshold is the weight above which two records potentially represent the same entity. The match threshold is the weight above which two records are considered to represent the same entity. Any records below the duplicate threshold are considered to represent completely separate and different entities.

Probabilities and Direct Weights

Optimum (or *ceiling*) matching weights can be assigned to field values using matching (m) and unmatching (u) probabilities or using agreement and disagreements weights in an equivalent way. Both types are based on a logarithmic function. Optimum agreement and disagreement weights are an equivalent logarithmic expression of the matching and unmatching probabilities, but for an end user, defining agreement and disagreement weight ranges is a more direct way to implement m-probabilities and u-probabilities.

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, the logarithmic formula using the m-probability and u-probability is negative, and a disagreement weight is subtracted from the composite match weight.

Agreement and Disagreement Weight Ranges

Like probabilities, the maximum agreement and minimum disagreement weights you define for each field allow you to specify 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 value assigned to a field is somewhere between the two numbers based on an underlying logarithmic formula. This provides a more convenient and intuitive representation of conditional probabilities.

Using the SSN and gender field example above, the SSN field is 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”.

How the Master Index Match Engine Works

The Master Index Match Engine compares records containing similar data types by calculating how closely certain fields in 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. You can also define custom comparison functions to plug in to the match engine. The formula used to determine the matching weight is based on either matching and unmatching probabilities or on agreement and disagreement weight ranges (described in [“Probabilities and Direct Weights” on page 9](#)).

The following topics provide additional information about how the Master Index Match Engine works:

- [“Master Index Match Engine Structure” on page 11](#)
- [“Master Index Match Engine Configuration Files” on page 11](#)
- [“Master Index Match Engine Matching Weight Formulation” on page 11](#)
- [“Master Index Match Engine Data Types” on page 12](#)
- [“The Master Index Match Engine and the Master Index Standardization Engine” on page 12](#)

Master Index Match Engine Structure

The Master Index Match Engine was designed to be very flexible and generic, allowing you to customize existing matching rules and to define additional rules using Java. The match engine framework allows you to create and plug in custom matching comparison functions, or *comparators*, to the match engine to enable matching against any type of data. The Master Index Match Engine framework includes two main modules. The real-time module stores the predefined and user-defined Java classes that define the matching comparator logic. The design-time module stores the configuration and validation classes for the comparators.

The Master Index Match Engine provides a wide variety of customizable comparators for you to choose from. You can also create comparators in the real-time module, and create new validation and configuration rules in the design-time module. The structure of the design-time module supports validations, weighting curves, and class dependencies. There is also an option that allows you load information from a data file and use that information to calculate a matching weight.

Master Index Match Engine Configuration Files

The Master Index Match Engine compares two records and returns a match weight indicating the likelihood of a match between the two records based on information provided in configuration files. In a master index application, the match engine is configured by these two files in the Match Engine node of the master index project: the matching configuration file (`matchConfigFile.cfg`) and the comparators list (`comparatorsList.xml`). The matching configuration file defines the configuration and parameters for the matching comparator functions and the comparators list defines each comparator available to the match engine.

Matching criteria and logic are defined in the match configuration file in the master index project (`matchConfigFile.cfg`). The data fields that are sent to the Master Index Match Engine for matching, known as the *match string*, are defined in the `MatchingConfig` section of `mefa.xml` in the master index project. The match engine configuration files define which matching rules to use to process each match field. The match engine provides a comprehensive set of comparator functions, and you can create custom comparators if needed.

Master Index Match Engine Matching Weight Formulation

The Master Index 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 Master Index 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.

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 weights are expressed as double values and can have up to 16 decimal points. When using agreement and disagreement weights, the Master Index Match Engine assigns a matching weight to each field that falls between the agreement and disagreement weights specified for the field. Thus, the maximum agreement weight between two records is the sum of the defined agreement weights for each field. The minimum disagreement weight is the sum of the defined disagreement weights for each field. For more information about weight calculation, see [“Determining the Weight Range” on page 62](#).

Master Index Match Engine Data Types

The Master Index Match Engine is built on a flexible framework that allows you to customize and create matching rules for various types of data. The match engine provides an extensive set of comparison functions for matching on various types of fields, such as numbers, dates, single characters, and so on. The match engine also provides more specialized comparison functions for searching on specific types of data, such as person names, address fields, social security numbers, genders. You can define custom comparison functions and custom standardization logic for different data types or variants on data types. These customizations are easily incorporated into a master index application, allowing you to completely customize the match and standardization process for your specific data format.

The Master Index Match Engine and the Master Index Standardization Engine

The Master Index Match Engine works with the Master Index Standardization Engine to provide an accurate comparison of two records. The standardization engine reads input data and determines how to parse, normalize, and phonetically encode the data in order to create a standard set of values to use for match comparison. The standardization engine can standardize free-form text fields, such as street address fields or business names, and separate them into their individual parts, such as house numbers, street names, and so on, allowing the match engine to generate a more accurate weight for free-form data.

Sun Master Index Standardization and Matching Process

In a default Sun Master Index implementation, the master index application uses the Master Index Match Engine and the Master Index Standardization Engine to cleanse data in real time. The standardization engine uses configurable pattern-matching logic to identify data and reformat it into a standardized form. The match engine uses a matching algorithm with a proven methodology to process and weight records in the master index database. By incorporating both standardization and matching capabilities, you can 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.

In a master index application, 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 these steps are complete, 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.

In a master index application, the standardization and matching process includes the following steps:

1. The master index application receives an incoming record.
2. The Master Index Standardization Engine standardizes the fields specified for parsing, normalization, and phonetic encoding. These fields are defined in `mefa.xml` and the rules for standardization are defined in the standardization engine 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 `master.xml`. 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 `mefa.xml`) and sends the string to the Master Index Match Engine.
5. The Master Index 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.

Master Index Match Engine Matching Configuration

The matching configuration files define how the Master Index Match Engine processes records to assign matching probability weights, allowing the master index application to identify matches, potential duplicates, and non-matches. The match engine includes two configurable files, the match configuration file and the comparators list. Together these files define additional logic for the Master Index Match Engine to use when determining the matching probability between two records.

The matching configuration is very flexible, allowing you to customize the matching logic according to the type of data being matched and for the record matching requirements of your business. In a master index application, the matching configuration files are stored in the master index project and are located in the Match Engine node of the project. The Master Index Standardization Engine typically standardizes the data prior to matching, so the match process is performed against the standardized data.

The matching configuration files for the Master Index 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 Master Index Match Engine Match Configuration File” on page 14](#)
- [“Master Index Match Engine Comparator Definition List” on page 21](#)

The Master Index 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 Master Index 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:

- “Master Index Match Engine Match Configuration File Format” on page 15
- “Master Index Match Engine Matching Comparison Functions At a Glance” on page 18
- “Master Index Match Engine Comparator Definition List” on page 21

These topics describe the format of the files so you can modify them directly. You can also modify the match configuration file using the Master Index Configuration Editor, which provides an easy, graphical way to configure matching rules.

Master Index 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. In a master index application, this file can be modified from the Matching tab of the Master Index Configuration Editor. For more information, see “Configuring the Comparison Functions for a Master Index Application” in *Sun Master Index Configuration Guide* [Configuring the Comparison Functions \(https://open-dm-mi.dev.java.net/docs/config/cnfg_index-eng-comp_t\)](https://open-dm-mi.dev.java.net/docs/config/cnfg_index-eng-comp_t)

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
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 Section

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.

For more information about probabilities and weights, see [“Probabilities and Direct Weights” on page 9](#).

Matching Rules Section

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. These are the rules you specify in the match string you define for a master index application. The syntax for this section is:
match-field size null-field function m-prob u-prob agreement disagreement
params data-sources

The following table describes each element in a match field row.

TABLE 1 Match Configuration File Columns

Column Number	Column Name	Description
1	match-field	A value that indicates to the Master Index Match Engine how each field should be weighted. Each field included in the match string (the MatchingConfig section of mefa.xml) 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”.

TABLE 1 Match Configuration File Columns (Continued)

Column Number	Column Name	Description
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"> ■ 0 - (zero) If one or both fields are empty, the weight used for the field is 0 (zero). ■ 1 - (one) If both fields are empty, the agreement weight is used; if only one field is empty, the disagreement weight is used. ■ a# - 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. ■ d# - 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. <p>Note – In the above descriptions, the agreement and disagreement weights are either specified in the file or calculated using a logarithmic formula based on the m and u-probabilities (depending on the probability type).</p>
4	function	The type of comparison to perform when weighting the field. For information about the available comparison functions, see Master Index Match Engine Comparison Functions and Options . An overview of the comparison functions is provided Table 2
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	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.

TABLE 1 Match Configuration File Columns (Continued)

Column Number	Column Name	Description
8	disagreement	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	params	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 Master Index Match Engine Comparison Functions and Options .
10	dataSources	The complete path to any data sources used by the comparison function specific in column 4. You can define as many data sources as there are data sources listed for the comparator in the comparators list file. The default comparators do not use data sources, but you can create a custom comparator that does.

Master Index Match Engine Matching Comparison Functions At a Glance

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 Master Index Match Engine should match the fields in the match string. The Master Index Match Engine provides several options to use with each function. You can also define custom comparison functions. For more information, see [Creating Custom Comparators for the Master Index Match Engine](#).

The following table summarizes each comparison function. A complete reference of the comparison functions and their parameters is included in [Master Index Match Engine Comparison Functions and Options](#).

Note – The names of these comparison functions are configurable. The following table lists their default names.

TABLE 2 Comparison Function Summary

Comparison Function	Name	Description
b1	Bigram Comparator	Compares two strings using an algorithm 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 Comparator	Compares two strings allowing for character transpositions. This function is similar to the standard Bigram Comparator (b1).
u	Advanced Jaro String Comparator	Compares two strings taking into account uncertainty factors, such as string length, transpositions, and characters in common. This function is based on the <i>Jaro</i> algorithm.
ua	Winkler-Jaro String Comparator	Compares two strings similar to the Advanced Jaro String Comparator (u), but increases the agreement weight if the initial characters of each string are exact matches. This function takes into account key punch and visual memory errors. It is based on the <i>Jaro</i> algorithm with variants of Winkler/Lynch and McLaughlin.
uf	Advanced Jaro Adjusted for First Names	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	Advanced Jaro Adjusted for Last Names	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	Advanced Jaro Adjusted for 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	Condensed String Comparator	Compares two strings similar to the Advanced Jaro String Comparator (u), but this function is 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 Advanced Jaro String Comparator, this function handles diacritical marks. This function also improves processing speed.
usu	Unicode String Comparator	Compares two strings similar to the Condensed String Comparator (us), but this function is based in Unicode to support multiple languages and alphabets. This comparator takes one parameter indicating the language to use.

TABLE 2 Comparison Function Summary *(Continued)*

Comparison Function	Name	Description
usus	Unicode AlphaNumeric Comparator	Compares two strings similar to the Unicode String Comparator, but this function is designed to match on unique identifiers such as national IDs. This comparator takes one parameter indicating the language to use plus any of the following parameters: <ul style="list-style-type: none"> ■ Field length ■ Character types ■ Invalid values
ujjs	Advanced Jaro AlphaNumeric Comparator	Compares two strings similar to the Advanced Jaro String Comparator, but this function is designed to match on unique identifiers such as national IDs. This comparator takes any of the following parameters: <ul style="list-style-type: none"> ■ Field length ■ Character types ■ Invalid values
c	Exact Character-to-Character Comparator	Compares string fields character by character. Each character must match in order for an agreement weight to be assigned.
nI	Integer Comparator	Compares integer 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.
nR	Real Number Comparator	Compares fields containing real numbers 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.
nS	Condensed AlphaNumeric SSN Comparator	Compares social security numbers or other unique identifiers, taking into account any of these parameters: <ul style="list-style-type: none"> ■ Field length ■ Character types ■ Invalid values

TABLE 2 Comparison Function Summary *(Continued)*

Comparison Function	Name	Description
dY	Date Comparator With Years as Units	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 With Months as Units	Compares the month and year using a relative distance as described above for the year comparison function (dY).
dD	Date Comparator With Days as Units	Compares the day, month, and year using a relative distance as described above for the year comparison function (dY).
dH	Date Comparator With Hours as Units	Compares the hour, day, month, and year using a relative distance as described above for the year comparison function (dY).
dm	Date Comparator With Minutes as Units	Compares the minute, hour, day, month, and year using a relative distance as described above for the year comparison function (dY).
ds	Date Comparator With Seconds as Units	Compares the second, minute, hour, day, month, and year using a relative distance as described above for the year comparison function (dY).
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.

Master Index Match Engine Comparator Definition List

The comparator definition list defines each comparator that is included in a master index application. If a comparator is not included in this list, it cannot be used in the application. If you define a comparator in this list that is not provided with the Master Index Match Engine, you need to define the logic of the new comparator in Java classes (for more information, see [“Creating Custom Comparators for the Master Index Match Engine”](#) on page 37).

Below is an excerpt from the default comparators list file that defines two numeric comparators, Real Number Comparator and Integer Comparator. Both comparators take two parameters, and are dependent on a second comparator class named CondensedStringComparator.

```

<comparator description="Numerics comparator">
  <className>NumericsComparator</className>
  <codes>
    <code description="Real Number Comparator" name="n[R, ]"/>
    <code description="Integer Comparator" name="nI" />
  </codes>
  <params>
    <param description="distance/string comparison option"
      name="switch" type="java.lang.String"/>
    <param description="Spectrum of comparison"
      name="range" type="java.lang.Integer|java.lang.Double"/>
  </params>
  <data-sources/>
  <dependency-classes>
    <dependency-class matchfield="CSC"
      name="com.sun.mdm.matcher.comparators.base.CondensedStringComparator"/>
  </dependency-classes>
  <curve-adjust status="false"/>
</comparator>

```

The comparators are defined in XML format. The following table lists and describes each element in the XML file.

TABLE 3 Comparator Definition List Elements

Element	Attribute	Description
group		An element that contains a list of comparators that all share the same Java package.
	description	A brief description of the comparator group.
	path	The Java package that contains the code that defines the comparators in the group.
comparator		A definition for one subgroup of comparators that are all based on the same Java class, have the same Java class dependencies, accept the same parameters and data sources, and have the same curve adjustment setting.
	description	A brief description of the comparator subgroup.
className		The name of the class that defines the logic for the comparators. The class must be contained in the package specified for the group element, as described above.
codes		A container element for a list of the comparators in the subgroup, with descriptions and processing codes of each comparator.

TABLE 3 Comparator Definition List Elements (Continued)

Element	Attribute	Description
code		A description and processing code for one comparator.
	description	A description of the comparator. The value you specify here appears in the comparator drop-down list on the Master Index Configuration Editor.
	name	A unique identifying name for the comparator. These are the comparator names used in the rules definitions in the match configuration file (<code>matchConfigFile.cfg</code>).
params		A container element for a list of static parameters for the subgroup of comparators. Parameters are optional.
param		One parameter definition for the comparators.
	description	A brief description of the parameter.
	name	A short name for the parameter.
	type	The Java data type of the values that can be specified for the parameter.
data-sources		A container element for a list of data files that contain additional information for the subgroup of comparators. For example, a comparator that generates weights based on the distance between postal codes might use lookup files containing information about the zip codes. Data sources are optional.
data-source		A definition for one data source. Currently, only file data sources are supported.
	description	A brief description of the data source.
	name	The complete path and filename of the data source.
	type	The type of data source being used. Currently, the only value you can specify is “java.io.File”.
dependency-classes		A container element that defines a list of Java classes on which the comparator class is dependent. The current comparator class inherits from the comparator classes you specify here as well as all the match fields (defined in <code>matchConfigFile.cfg</code>) that use that comparator.
dependency-class		A definition for one comparator class, called a <i>dependency comparator</i> , on which the current comparator class is dependent.
	matchField	The name of the dependency comparator’s match field.
	name	The name of the dependency comparator class.

TABLE 3 Comparator Definition List Elements (Continued)

Element	Attribute	Description
curve-adjust		An indicator of whether to apply special adjustments to the weighting curve. The curve adjustment is defined for each comparator individually in a Java class named <i>comparator_nameCurveAdjustor</i> .
	status	The status of the curve adjustor. Specify true to use the curve adjustor; specify false to disable the curve adjustor.

Master Index Match Engine Comparison Functions

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 Master Index Match Engine matches records. The comparators themselves are highly configurable and can be configured to assign differing weights or handle null values. Several comparators accept parameters that further fine-tune the matching process.

The Master Index Match Engine provides a comprehensive group of match comparison functions to enable matching on a wide variety of data. While you should be able to configure any of the default comparison functions to accurately match your data, you can create new comparison functions and integrate them into a master index application. For more information, see [“Creating Custom Comparators for the Master Index Match Engine” on page 37](#).

The following topics describe each of the default comparison functions:

- [“Bigram Comparators” on page 25](#)
- [“Uncertainty String Comparators” on page 25](#)
- [“Exact Character-to-Character Comparator \(c\)” on page 30](#)
- [“Numeric Comparators” on page 31](#)
- [“Condensed AlphaNumeric SSN Comparator \(nS\)” on page 32](#)
- [“Date Comparators” on page 33](#)
- [“Prorated Comparator \(p\)” on page 35](#)

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, parameters define how null fields are handled and what the actual agreement and disagreement weights are.

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

Bigram Comparators

The Master Index 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 Comparator (b1)

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

Advanced Bigram Comparator (b2)

The Advanced Bigram Comparator 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 Master Index Match Engine provides several uncertainty comparison functions for comparing string fields. Most uncertainty comparison functions are generic, but some comparison functions are designed for specific types of information (first name, last name, house number, and national identifiers).

The uncertainty functions include the following:

- “Advanced Jaro String Comparator (u)” on page 26
- “Winkler-Jaro String Comparator (ua)” on page 26
- “Condensed String Comparator (us)” on page 27
- “Advanced Jaro Adjusted for First Names (uf)” on page 27
- “Advanced Jaro Adjusted for Last Names (ul)” on page 27
- “Advanced Jaro Adjusted for House Numbers (un)” on page 27
- “Advanced Jaro AlphaNumeric Comparator (ujs)” on page 27

- “Unicode String Comparator (usu)” on page 28
- “Unicode AlphaNumeric Comparator (usus)” on page 29

Advanced Jaro String Comparator (u)

The Advanced Jaro String Comparator is the standard uncertainty comparison function for processing string fields. This 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.

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.

Winkler-Jaro String Comparator (ua)

The Winkler-Jaro String Comparator 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, inserting “1” the number instead of “l” the letter) or keypunch errors (for example, typing “S” instead of “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.

Condensed String Comparator (us)

The Condensed String Comparator is a custom version of a generic string comparison function. It is similar to the Advanced Jaro String Comparator, **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 Advanced Jaro String Comparator, the Condensed String Comparator 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.

Advanced Jaro Adjusted for First Names (uf)

The Advanced Jaro Adjusted for First Names comparator is designed specifically for matching on first name fields, and is based on the Condensed String Comparator, **us**. This comparison function analyzes the string and then adjusts the weight based on statistical data. This comparison function takes no parameters.

Advanced Jaro Adjusted for Last Names (ul)

The Advanced Jaro Adjusted for Last Names comparator is designed specifically for matching on last name fields, and is based on the Condensed String Comparator, **us**. This comparison function analyzes the string and then adjusts the weight based on statistical data. This comparison function takes no parameters.

Advanced Jaro Adjusted for House Numbers (un)

The Advanced Jaro Adjusted for House Numbers comparator is designed specifically for matching on house numbers, and is based on the Condensed String Comparator, **us**. This comparison function analyzes the string and then adjusts the weight based on statistical data. This comparison function takes no parameters.

Advanced Jaro AlphaNumeric Comparator (ujs)

The Advanced Jaro AlphaNumeric Comparator is a custom version of a generic string comparison function. It is based on the Advanced Jaro String Comparator, **u**, but is designed specifically for matching on national identifier, such as social security numbers. This function takes into account such uncertainty factors as string length, transpositions, key punch errors,

and visual memory errors. It can also take into consideration field length, allowed character types, and invalid values. This comparison function takes the parameters described in the following table.

TABLE 4 ujs Comparison Function Parameters

Parameter	Description
ssnLength	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 field size in the matching configuration file (for more information, see “Matching Rules Section” on page 16).
recType	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.
ssnList	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.

Unicode String Comparator (usu)

The Unicode String Comparator is a custom version of a generic string comparison function. It is similar to the Condensed String Comparator, **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 the following table.

TABLE 5 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>da - Danish</p> <p>sv - Swedish</p> <p>nb - Norwegian Bokmål</p> <p>nn - Norwegian Nynorsk</p> <p>nl - Dutch</p> <p>es - Spanish</p> <p>fr - French</p> <p>en - English</p> <p>it - Italian</p> <p>de - German</p>

Unicode AlphaNumeric Comparator (usus)

This comparison function is a custom version of a generic string comparison function. It is similar to the Unicode String Comparator, but it is also similar to the Advanced Jaro AlphaNumeric Comparator in that it is designed to work on national identifiers like social security numbers. This locale-oriented comparator recognizes the nuances of each language and supports the complexities and subtleties of each. This function takes into account such uncertainty factors as string length, transpositions, key punch errors, and visual memory errors. It can also take into consideration field length, allowed character types, and invalid values. This comparison function takes the parameters described in the following table.

TABLE 6 usus Comparison Function Parameters

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>da - Danish</p> <p>sv - Swedish</p> <p>nb - Norwegian Bokmål</p> <p>nn - Norwegian Nynorsk</p> <p>nl - Dutch</p> <p>es - Spanish</p> <p>fr - French</p> <p>en - English</p> <p>it - Italian</p> <p>de - German</p>
fixed-length	<p>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 “Matching Rules Section” on page 16).</p>
character-type	<p>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.</p>
invalid-characters	<p>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.</p>

Exact Character-to-Character Comparator (c)

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

Numeric Comparators

The Master Index Match Engine provides two comparison functions for matching on numeric fields:

- “Integer Comparator (nI)” on page 31
- “Real Number Comparator (nR)” on page 32

The Integer Comparator and Real Number Comparator can perform either 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 1 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.

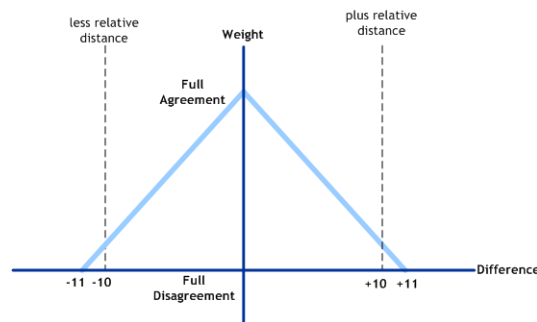


FIGURE 1 Numeric Relative Distance Comparison

Integer Comparator (nI)

The Integer Comparator matches specifically on integers using the logic describe above. It accepts the parameters listed in the following table.

TABLE 7 nI Comparison Function Parameters

Parameter	Description
switch	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.
range	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).

Real Number Comparator (nR)

The Real Number Comparator function matches specifically on real numbers based on the logic described above. It accepts the parameters listed in the following table.

TABLE 8 nR Comparison Function Parameters

Parameter	Description
switch	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.
range	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).

Condensed AlphaNumeric SSN Comparator (nS)

The Condensed AlphaNumeric SSN Comparator is designed specifically for matching on numeric strings and is very useful for matching social security numbers or other unique identifiers. This comparison function can compare either alphanumeric values or numeric values, and takes into account such uncertainty factors as string length, transpositions, key punch errors, and visual memory errors. It can also take into consideration field length, allowed character types, and invalid values. It accepts the parameters listed in [Table 9](#).

TABLE 9 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 “Matching Rules Section” on page 16).
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.
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 Master Index 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 With Years as Units \(dY\)” on page 34](#)
- [“Date Comparator With Months as Units \(dM\)” on page 35](#)
- [“Date Comparator With Days as Units \(dD\)” on page 35](#)
- [“Date Comparator With Hours as Units \(dH\)” on page 35](#)
- [“Date Comparator With Minutes as Units \(dm\)” on page 35](#)
- [“Date Comparator With Seconds as Units \(ds\)” on page 35](#)

As with the numeric comparison functions, the date comparison functions can use either a direct string comparison or a relative distance calculation (see [“Numeric Comparators” on page 31](#)). 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 2](#) 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.

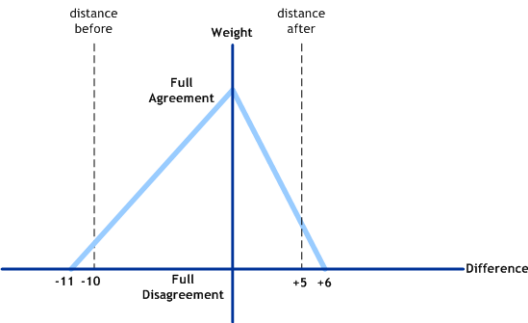


FIGURE 2 Date Relative Distance Comparison

The date comparison functions take the parameters listed in [Table 10](#).

TABLE 10 Date Comparison Function Parameters

Parameter	Description
switch	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.
llimit	The number of units prior to the reference date/time for which two date fields can still be considered a match.
ulimit	The number of units following the reference date/time for which two date fields can still be considered a match.

Date Comparator With Years as Units (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 With Months as Units (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 With Days as Units (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 With Hours as Units (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 With Minutes as Units (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 With Seconds as Units (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 The Prorated Comparator 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 3 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.

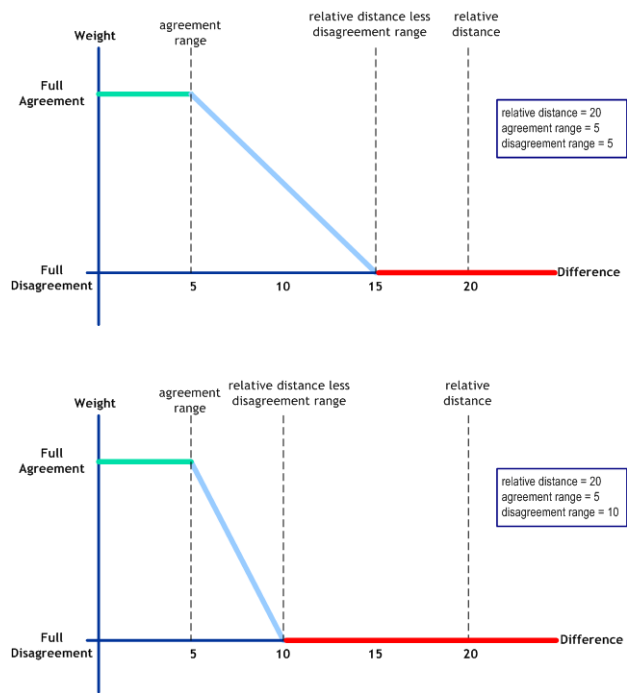


FIGURE 3 Prorated Linear Adjustment Comparison

The prorated comparison functions takes the parameters listed in Table 11.

TABLE 11 Prorated Comparison Function Parameters

Parameter	Description
range	The greatest difference between two numbers at which they can still be considered a match or partial match.
tolerance1	The greatest difference between two numbers at which they are considered a full match. This number must be less than the relative distance.

TABLE 11 Prorated Comparison Function Parameters (Continued)

Parameter	Description
tolerance2	<p>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.</p> <p>The weighting scale decreases in size as the value of the full-disagreement parameter increases (see diagram).</p>

Creating Custom Comparators for the Master Index Match Engine

The Master Index Match Engine provides a variety of configurable matching comparators for you to process and match you data. However, if none of the existing comparators meet your requirements, the flexible framework of the Master Index Match Engine allows you to create custom comparators to plug in to master index applications. The comparators are flexible components that can be modified and tailored without requiring any changes to the framework.

The following topics provide information and instructions about creating your own custom comparators:

- [“Custom Comparator Overview” on page 37](#)
- [“About the Comparator Package” on page 38](#)
- [“Defining Custom Comparators” on page 39](#)

Custom Comparator Overview

Creating a custom matching comparator for the Master Index Match Engine requires coding the processing and validation logic for the comparator in Java. The Master Index Match Engine provides the interfaces and supporting Java classes you need to implement in order to incorporate the comparators into a master index application.

The Master Index Match Engine framework consists of two modules. The real-time module stores the basic logic for the matching comparators. The design-time module stores all of the configuration logic for the comparators, including parameter validations, data source definitions, and curve adjustment logic. The two pieces are pulled together by the configuration in the comparators list file (`comparatorsList.xml`). For each custom comparator package you create, you need to create a comparators list file.

You can define the following information in the comparators list for each comparator you create.

- A code that is used to reference the comparator in the match configuration file (`matchConfigFile.cfg`).
- The class that defines the comparator logic.
- Parameters for the comparator. Parameter values are entered in the match configuration file for any entries that reference the comparator.
- Any classes from which the comparator class inherits.
- Data sources that provide additional information to the comparator during the match process.
- Whether to use curve adjustment logic for the comparator.

After you create the package, you can import the custom comparators into NetBeans using the easy import function of Sun Master Index. When you import the files, Sun Master Index automatically validates the files and merges the comparators list information into the comparators list for the application. You can then add and configure entries in `matchConfigFile.cfg` that reference the comparator, which makes the comparator available to be used in the match string.

About the Comparator Package

After you register your custom comparators and you create and compile the comparators and any configuration classes, you need to package the files in a ZIP file so they are available for import into NetBeans. For optimal usage, it is best to package all similar comparators in a unique ZIP file. You can create single packages for each comparator, or combine them into one package.

The ZIP file includes the following:

- The comparator Java classes
- The comparators list file (`comparatorsList.xml`)
- Any parameter validation classes (only if the comparators take parameters)
- Any data source loading or validation classes (only if the comparators use external data files)
- Any curve adjustment classes (only if the comparators use curve adjustment for weight calculation)

For the ZIP file to have the correct structure, the `comparatorsList.xml` file should be at the same level as the `com` folder that contains the Java classes. The following figure shows a sample ZIP file for custom comparators.

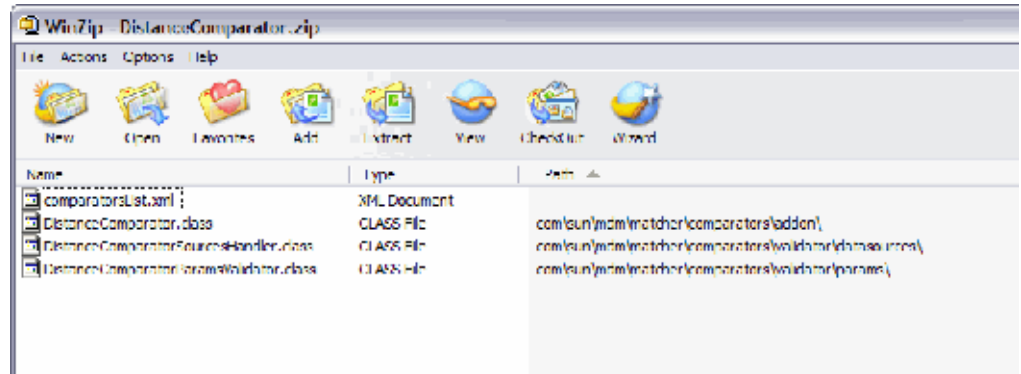


FIGURE 4 Sample Custom Comparator ZIP Package

Defining Custom Comparators

The following topics provide instructions for each step of creating custom comparators. You might need to create multiple Java files and Java packages for the comparator, depending on the validations, data sources, dependency classes, and curve adjustments you use. Create them in the same directory structure because you will need to package them up into a ZIP file when you are through.

- “Before You Begin” on page 39
- “Step 1: Create the Custom Comparator Java Class” on page 40
- “Step 2: Register the Comparator in the Comparators List” on page 43
- “Step 3: Define Parameter Validations (Optional)” on page 45
- “Step 4: Define Data Source Handling (Optional)” on page 46
- “Step 5: Define Curve Adjustment or Linear Fitting (Optional)” on page 50
- “Step 6: Compile and Package the Comparator” on page 51
- “Step 7: Import the Comparator Package Into Sun Master Index” on page 52
- “Step 8: Configure the Comparator in the Match Configuration File” on page 52

Before You Begin

Before you create your custom comparators, take into account the following requirements for the comparators.

- Determine how many comparators you need to create and whether each will require a different Java class or some can use the same Java class.
- Determine what parameters, if any, you need to define for each comparator.
- Determine what validations, if any, need to be created.
- Determine whether you need to use a data source.

- Decide if the comparators you create will have a dependency on any other comparator classes.
- Decide whether you will use curve adjustment, linear fitting, or neither.

Step 1: Create the Custom Comparator Java Class

The first step to creating custom comparators is defining the matching logic in custom comparator Java classes that are stored in the real-time module of the Master Index Match Engine. Follow these guidelines when creating the class:

- Create a working directory that will contain all the Java packages and the comparators list file for the new comparators.
- The Java classes need to implement `com.sun.mdm.matcher.comparators.MatchComparator.java` interface, located in `Matcher.jar`. This class includes the methods described below.

Once you create the Java classes, continue to [“Step 2: Register the Comparator in the Comparators List” on page 43](#).

initialize

Description

The `initialize` method initializes the values for the parameters, data sources, and dependency class used for each custom comparator. It provides the necessary information to access the comparator's configuration in the match configuration file and the comparators list file.

Syntax

```
void initialize(Map<String, Map> params, Map<String, Map> dataSources,  
Map<String, Map> dependClassList)
```

Parameters

Parameter	Type	Description
params	Map	A mapping of all the parameters associated with a match field in <code>matchConfigFile.cfg</code> .
dataSources	Map	A mapping of all the data sources associated with a match field in <code>matchConfigFile.cfg</code> .

Parameter	Type	Description
dependClassList	Map	A mapping of all the dependency classes associated with a match field in <code>matchConfigFile.cfg</code> .

Return Value

None.

Throws

None.

compareFields

Description

The `compareFields` method contains all the comparison logic needed to compare two field values and calculate a matching weight that shows how similar the values are.

Syntax

```
double compareFields(String recordA, String recordB, Map context)
```

Parameters

Parameter	Type	Description
recordA	String	A field value from the record against which the reference record is being compared.
recordB	String	A field value from the reference record.
context	Map	A set of arguments passed to the comparator.

Return Value

A number between zero and one that indicates how closely two field values match.

Throws

`MatchComparatorException`

setRTParameters

Description

The `setRTParameters` method sets the runtime parameters for the comparator, providing the ability to customize every call to the parameter.

Syntax

```
void setRTParameters(String key, String value)
```

Parameters

Parameter	Type	Description
key	string	The key to map the parameter value.
value	string	The value of the parameter.

Return Value

None.

Throws

None.

stop

Description

The `stop` method closes any related connections to the data sources used by the comparator.

Syntax

```
void stop()
```

Parameters

None.

Return Value

None.

Throws

None.

Step 2: Register the Comparator in the Comparators List

In order to include new comparators in a master index application, you need to create a comparators list file defining the configuration of the comparators. When you import the comparator package into the master index application, this file is read and the entries are added to the comparators list for the project.

Below is a sample comparators list file. Note that the first comparator includes all possible configurations (parameters, dependency classes, data sources, and curve adjust). Most comparators will not be that complex. The second comparator class defines two comparators, Approx and Adjust.

```
<?xml version="1.0" encoding="UTF-8"?>
<comparators-list xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
xsi:noNamespaceSchemaLocation="comparatorsList.xsd">
  <group description="New group of comparators"
    path="com.mycomparators.matchcomparators">
    <comparator description="New Exact Comparator">
      <className>NewExactComparator</className>
      <codes>
        <code description="New Exact Comparator" name="Exact" />
      </codes>
      <params>
        <param description="Fixed length" name="length"
          type="java.lang.Integer" />
        <param description="Data type" name="dataType"
          type="java.lang.String" />
      </params>
      <data-sources>
        <datasource description="Serial numbers" type="java.io.File" />
      </data-sources>
      <dependency-classes>
        <dependency-class matchfield="Serial"
          name="com.genericcomparaotrs.StringComparator" />
      </dependency-classes>
      <curve-adjust status="true" />
    </comparator>
    <comparator description="New Approximate Comparator">
      <className>NewApproxComparator</className>
      <codes>
        <code description="New approximate comparator" name="Approx" />
```

```
        <code description="New adjustable comparator" name="Adjust" />
      </codes>
    </comparator>
  </group>
</comparators-list>
```

▼ To Register the Comparators

- 1 Complete [“Step 1: Create the Custom Comparator Java Class” on page 40.](#)
- 2 In the same folder where you created the custom Java class package, create a new file named `comparatorsList.xml`.

Tip – The comparators list file needs to be in the same working directory you created for the custom comparator Java classes.

- 3 Add the following header information to the file. You can copy this from the `comparatorList.xml` file in a master index application.

```
<?xml version="1.0" encoding="UTF-8"?>
<comparators-list xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
  xsi:noNamespaceSchemaLocation="comparatorsList.xsd">
  ...
</comparators-list>
```

- 4 Define the following properties, using the XML structure described in [“Master Index Match Engine Comparator Definition List” on page 21](#). Use the sample above as an example.
 - The group description and Java package for the group.
 - A description for each comparator.
 - The Java class name for each comparator or comparator subgroup.
 - The unique identifying name for each comparator.
 - A list of static parameters for each comparator or comparator subgroup (optional). If you define parameters, you must also perform the steps under [“Step 3: Define Parameter Validations \(Optional\)” on page 45](#).
 - A list of data sources for each comparator or comparator subgroup (optional). If you define data sources, you must also perform the steps under [“Step 4: Define Data Source Handling \(Optional\)” on page 46](#).
 - A list of dependency classes for each comparator or comparator subgroup (optional).
 - Whether to use curve adjustment for each comparator or comparator subgroup (optional). If you set curve adjustment to true, you must perform the steps under [“Step 5: Define Curve Adjustment or Linear Fitting \(Optional\)” on page 50](#).

- 5 Continue to [“Step 3: Define Parameter Validations \(Optional\)” on page 45](#)

Step 3: Define Parameter Validations (Optional)

If your custom comparators take parameters, you should create a Java class that validates the parameter properties. You need to perform this step if you defined parameters for the comparator in `comparatorsList.xml`. You do not need to create this file in the same package as the Java comparator class, but for packaging purposes, create it in the same working folder.

▼ To Define Parameter Validations

- 1 Complete [“Step 2: Register the Comparator in the Comparators List” on page 43](#).
- 2 Create a Java class named the same name as the Java class that defines the comparator with `ParamsValidator` appended.
For example, if the comparator is defined by a class named `ExactComparator`, the parameter validation class would be `ExactComparatorParamsValidator`.
- 3 In this class, implement `com.sun.mdm.matcher.comparators.validator.ParametersValidator`.
The method contained in this class is described below.
- 4 Continue to [“Step 4: Define Data Source Handling \(Optional\)” on page 46](#).

validateComparatorsParameters

Description

The `ParametersValidator` class contains one method, `validateComparatorsParameters`, that allows you to validate parameter types, ranges, and other properties. For logging purposes, you can use `net.java.hulp.i18n`, which is used within `matcher.jar`, or you can use your own logger.

Syntax

```
void validateComparatorsParameters(Map<String, Object> params)
```

Parameters

Parameter	Type	Description
params	Map	A list of parameters to validate.

Return Value

None.

Throws

MatcherException

Step 4: Define Data Source Handling (Optional)

If your custom comparators use external data sources to provide additional information for matching weight calculations, you need to create a Java class that lets you load the file to memory or have real-time access to the data file content. You can also define validations to perform. You do not need to create this file in the same package as the Java comparator class, but for packaging purposes, create it in the same working folder.

You need to perform this step if you defined lines similar to the following in `comparatorsList.xml`:

```
<data-sources>
  <datasource description="Serial numbers" type="java.io.File" />
</data-sources>
```

▼ To Define Data Source Handling

- 1 Complete [“Step 3: Define Parameter Validations \(Optional\)” on page 45.](#)
- 2 Create a Java class named the same name as the Java class that defines the comparator with **“SourcesHandler”** appended.

For example, if the comparator is defined by a class named `ExactComparator`, the parameter validation class would be `ExactComparatorSourcesHandler`.
- 3 In this class, implement
`com.sun.mdm.matcher.comparators.validator.DataSourcesHandler`.
The method in this class is described below.
- 4 Continue to [“Step 5: Define Curve Adjustment or Linear Fitting \(Optional\)” on page 50.](#)

handleComparatorsDataSources

Description

The DataSourcesHandler class contains one method, handleComparatorsDataSources, that allows you to define properties for the data source. This method takes one parameter that is a DataSourcesProperties object. This class and its methods are described in [“DataSourcesProperties Class” on page 47](#).

Syntax

Object handleComparatorsDataSources(DataSourcesProperties dataSources)

Parameters

Parameter	Type	Description
dataSources	DataSourceProperties	A list of properties for the data handler (see “DataSourcesProperties Class” on page 47).

Return Value

Object

Throws

MatcherException

IOException

DataSourcesProperties Class

The DataSourcesProcerties interface is used as a parameter to the handleComparatorsDataSources described in [“Step 4: Define Data Source Handling \(Optional\)” on page 46](#). The methods in the class are listed and described below.

getDataSourcesList

Description

The getDataSourcesList returns the comparator's list of associated data source paths.

Syntax

List getDataSourcesList(String codeName)

Parameters

Parameter	Type	Description
codeName	string	The name of the comparator. The name is defined in <code>comparatorsList.xml</code> in the <code>name</code> attribute of the <code>code</code> element. In the example below, the comparator's code name is "Exact". <code><code description="New exact comparator" name="Exact" /></code>

Return Value

A list of paths and filenames as specified in `comparatorsList.xml`.

Throws

None.

isDataSourceLoaded

Description

The `isDataSourceLoaded` method checks whether a specific file has already been loaded or opened.

Syntax

`boolean isDataSourceLoaded(String sourcePath)`

Parameters

Parameter	Type	Description
sourcePath	string	The path and filename of the file to check.

Return Value

A boolean indicator of whether the specified file has already been loaded or opened.

Throws

None.

setDataSourceLoaded

Description

The `setDataSourceLoaded` method sets the loading status of a data source.

Syntax

```
void setDataSourceLoaded(String sourcePath, boolean status)
```

Parameters

Parameter	Type	Description
sourcePath	string	The path and filename of the file.
status	boolean	The load status of the file. Specify true if the file is loaded; otherwise specify false .

Return Value

None.

Throws

None.

getDataSourceObject

Description

The `getDataSourceObject` method returns the file located at the specified source path.

Syntax

```
Object getDataSourceObject(String sourcePath)
```

Parameters

Parameter	Type	Description
sourcePath	string	The path and filename of the file you want to load.

Return Value

An object containing the data source information.

Throws

None.

Step 5: Define Curve Adjustment or Linear Fitting (Optional)

If your custom comparators use curve adjustment or linear fitting to adjust matching weight calculations, you need to create a Java class that defines the curve. You do not need to create this file in the same package as the Java comparator class, but for packaging purposes, create it in the same working folder.

You need to perform this step if you defined the following line in `comparatorsList.xml` for the comparator:

```
<curve-adjust status="true" />
```

▼ To Define Curve Adjustment or Linear Fitting

- 1 Complete [“Step 4: Define Data Source Handling \(Optional\)” on page 46.](#)
- 2 Create a Java class named the same name as the Java class that defines the comparator with **“CurveAdjustor”** appended.

For example, if the comparator is defined by a class named `ExactComparator`, the parameter validation class would be `ExactComparatorCurveAdjustor`.
- 3 In this class, implement `com.sun.mdm.matcher.configurator.CurveAdjustor`.
The method in this class is described below.
- 4 Continue to [“Step 6: Compile and Package the Comparator” on page 51.](#)

processCurveAdjustment

Description

The `processCurveAdjustment` method provides handling for curve adjustment within a specific match comparator.

Syntax

```
double[] processCurveAdjustment(String compar, double[] cap)
```

Parameters

Parameter	Type	Description
compar	string	The name of the comparator, as defined in the name attribute of the code element for the comparator.
cap	double[]	An array of values that define the curve adjustment.

Return Value

An array of curve adjustment values.

Throws

`MatcherException`

Step 6: Compile and Package the Comparator

Before you perform these steps, make sure you have completed [“Step 1: Create the Custom Comparator Java Class”](#) on page 40 through [“Step 5: Define Curve Adjustment or Linear Fitting \(Optional\)”](#) on page 50.

When you are finished defining all the Java classes for the comparators and have registered each comparator in your comparators list file, you can compile the Java code and package the files into a ZIP file that you can then import into a master index application. Compile the classes using the compiler of your choice.

To package the files, create a temporary directory and copy the comparators list file to the directory. Copy all the class folders and files to the same directory. The top level of the temporary directory should include `comparatorsList.xml` and a `com` folder (which contains all the Java classes). Create a ZIP file of the directory. For more information about the ZIP package, see [“About the Comparator Package”](#) on page 38.

After you compile and package the comparator, continue to [“Step 7: Import the Comparator Package Into Sun Master Index” on page 52.](#)

Step 7: Import the Comparator Package Into Sun Master Index

You need to import the your new comparators into NetBeans to make them available to all master index applications or only the current application.

▼ To Import a Comparison Function

- 1 Launch NetBeans, and open the master index project that will use the new comparators.
- 2 In the Projects window, expand the main master index project.
- 3 Right-click Match Engine, and select Import Comparator Plug-in.
- 4 In the dialog box that appears, navigate to the location of the plug-in ZIP file.
- 5 Select the file containing the plug-in, and then click Open.
- 6 Do one of the following:
 - To import the plug-in and make it available to all future master index application, click Yes.
 - To import the plug-in and make it only available to the current master index application, click No.

The contents of the ZIP file are imported into the Match Engine node and the new comparators are added to the list of comparator definitions in `comparatorsList.xml`.

- 7 In the Match Engine node, navigate to the `/lib` folder that was added and verify that all of the required files are there.
- 8 Open `comparatorsList.xml` and verify the new comparator definitions are included.

Step 8: Configure the Comparator in the Match Configuration File

After you import custom comparators, you need to add them to the match configuration file (`matchConfigFile.cfg`) and define the matching configuration. This makes the comparator

available for use in the master index match string. For information about this file, see [“The Master Index Match Engine Match Configuration File” on page 14](#). For instructions on modifying the file, see [“Configuring the Comparison Functions for a Master Index Application” in *Sun Master Index Configuration Guide*](#).

Master Index Match Engine Configuration for Common Data Types

The Master Index Match Engine can match on any type of data. Common data types for matching include person names, addresses, and business names. Configuring the match engine for matching on these data types in a master index application includes modifying the match configuration file (`matchConfigFile.cfg`) and `mefa.xml`.

The Master Index Match String

The data string that is passed to the Master Index Match Engine for match processing is called the *match string*. For a master index application, the match string is defined in the `MatchingConfig` section of `mefa.xml`. The match and standardization 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 match engine from the records in the candidate selection pool. Finally, the Master Index Match Engine configuration files define how the match string is processed.

The Master Index 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. For more information about match types, see [“Master Index Match Engine Match Types” on page 55](#)).

Master Index Match Engine Match String Fields

In a master index application, the match string processed by the Master Index Match Engine is defined by the match fields specified in `mefa.xml`, and the logic for how the fields are matched is defined in the match configuration file (`matchConfigFile.cfg`). The match engine can process any combination of fields you specify for matching using the predefined comparators or any new comparators you define. Not all fields in a record need to be processed by the Master Index Match Engine. Before you define the match string, analyze your data to determine the fields that are most likely to indicate a match or non-match between two records.

The following topics provide additional information about the match string for different data types:

- [“Person Data Match String Fields” on page 54](#)
- [“Address Data Match String Fields” on page 54](#)
- [“Business Name Match String Fields” on page 55](#)

Person Data Match String Fields

By default, the match configuration file (`matchConfigFile.cfg`) includes rows specifically for matching on first name, last name, social security numbers, and dates (such as a date of birth). It also includes a row for matching a single character with logic specialized for 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. When matching on person names, determine whether you want to use the original field values, the normalized field values, or the phonetic values. The match engine can handle any of these types of fields, but the best comparator for each type might be different. Also determine how much weight you want to give each field type and configure the match configuration file accordingly.

Address Data Match String Fields

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 configure for address fields. You can use any of the existing rows for matching or you can add rows for the fields you want to match. If you specify an “Address” match type for any field in the Master Index Wizard, the default fields that store the parsed data are automatically added to the match string in `mefa.xml`. These fields include the house number, street direction, street type, and street name. You can remove any of these fields from the match string.

When matching on address fields, determine whether you want to use the original field values, the standardized field values, or the phonetic values. The match engine can handle any of these

types of fields, but the best comparator for each type might be different. Also determine how much weight you want to give each field type and configure the match configuration file accordingly.

Business Name Match String Fields

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 customize to use with business name fields. You can use any of the existing rows for matching or you can add rows for the fields you want to match. 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 `mefa.xml`, including the name, organization type, association type, sector, industry, and URL. You can remove any of these fields from the match string.

When matching on business name fields, determine whether you want to use the original field values, the standardized field values, or the phonetic values. The match engine can handle any of these types of fields, but the best comparator for each type might be different. Also determine how much weight you want to give each field type and configure the match configuration file accordingly.

Master Index Match Engine Match Types

The default match configuration file, `matchConfigFile.cfg`, defines several rules that you can customize for the type of data being processed. Each rule is identified by a *match type* in the first column of each row. This value identifies the type of matching to perform to the match engine. In a master index application, the match type is entered for each field in the match string section of `mefa.xml`.

The match configuration Master Index Match Enginefile 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 62](#) and [“Master Index Match Engine Comparison Functions” on page 24](#).

The following four tables list match types that are typically used in processing different data types, including:

- [Person Data Match Types](#)
- [Address Match Types](#)
- [Business Name Match Types](#)
- [Miscellaneous Match Types](#)

The following match types are designed for matching on person data.

TABLE 12 Person Data Match Types

This indicator ...	processes this data type ...
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.
SSN	A field containing a social security number.
Gender	A field containing a gender code.

The following match types are designed for matching on address data.

TABLE 13 Address Match Types

This indicator ...	processes this data type ...
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.

The following match types are designed for matching on business names.

TABLE 14 Business Name Match Types

This match type ...	processes this data type ...
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.
LocationTypeKeyword	The parsed location 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.

Miscellaneous match types provide additional logic for matching on a variety of data types, such as date, numeric, string, and character fields.

TABLE 15 Miscellaneous Match Types

This indicator ...	processes this data type ...
Date	The year of a date field.
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.
String	A generic string field.
Unistring	A generic Unicode string field.
Integer	A field containing integers.
Real	A field containing real numbers.
Char	A field containing a single character.
pro	Any field on which you want the Master Index Match Engine to use prorated weights.
Exac	Any field you want the Master Index Match Engine to match character for character.
CSC	A generic string.
DOB	A date of birth in string rather than date format.

Configuring the Match String for a Master Index Application

The *MatchingConfig* section of *mefa.xml* determines which fields are passed to the Master Index Match Engine for matching (the match string). The match types specified in this section help the match engine determine the algorithm and custom logic to use for matching on each field.

If you are matching on fields parsed from a free-form text field, define each individual parsed field you want to use for matching in the Master Index Wizard or Configuration Editor. The match types you can use for each field in this section are defined in the first column of the match configuration file (*matchConfigFile.cfg*). Make sure the match type you specify has the correct matching logic defined in the match configuration file. See [“Master Index Match Engine Match Types” on page 55](#) for more information.

The following topics provide more information about matching on different types of data:

- [“Configuring the Match String for Person Data” on page 58](#)
- [“Configuring the Match String for Address Data” on page 59](#)
- [“Configuring the Match String for Business Names” on page 60](#)

Configuring the Match String for Person Data

When matching on person data, you can include any field stored in the database for matching. To configure the match string, follow the instructions under [“Defining the Master Index Match String” in *Sun Master Index Configuration Guide*](#). For the Master Index Match Engine, each data type has a different match type (specified by the *match-type* element in the matching configuration file). The FirstName, LastName, SSN, Gender, and DOB 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 [“Master Index Match Engine Match Types” on page 55](#).

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>

```

Configuring the Match String for Address Data

For matching on street address fields, make sure the match string you specify in the MatchingConfig section of mefa.xml contains all or a subset of the fields that contain the standardized data (the original text in street address fields is 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” in *Sun Master Index Configuration Guide*](#). For the Master Index Match Engine, each component of a street address has a different match type (specified by the *match-type* element in the matching configuration file). 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 [“Master Index Match Engine Match Types” on page 55](#).

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-columns>
</match-system-object>

```

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

Configuring the Match String for Business Names

For matching on business name fields, make sure the match string you specify in the MatchingConfig section of mefa.xml 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” in *Sun Master Index Configuration Guide*](#). For the Master Index Match Engine, each data type has a different match type (specified by the *match-type* element of the matching configuration file). 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 [“Master Index Match Engine Match Types” on page 55](#).

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>
  <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 Master Index

Each Sun Master Index 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 61](#)
- [“Customizing the Match Configuration and Thresholds” on page 61](#)

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 Master Index to analyze data, review the information in [Analyzing and Cleansing Data for Sun Master Index](#) and [Loading the Initial Data Set for a Sun Master Index](#) before you extract the legacy data.

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 61](#)
- [“Customizing the Match Configuration” on page 62](#)
- [“Determining the Weight Thresholds” on page 64](#)

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 [“Master Index Match Engine Matching Weight Formulation” on page 11](#) 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. Using the data analysis tool provided can help you determine the match and duplicate thresholds.

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. The following table 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 16 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
Maximum Match Weight	38	
Minimum Match Weight		-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:

$$\text{LOG2}((1-m_prob)/(1-u_prob))$$

The following table illustrates 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 17 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

TABLE 17 Sample m-probabilities and u-probabilities (Continued)

Field Name	m-probability	u-probability	Max Agreement Weight	Min Disagreement Weight
Maximum Match Weight			38	
Minimum Match Weight				-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 [“Master Index Match Engine Comparison Functions” on page 24](#). 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 Master Index Match Engine. Two parameters in master.xml 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 5 illustrates the match and duplicate thresholds in comparison to total composite match weights.

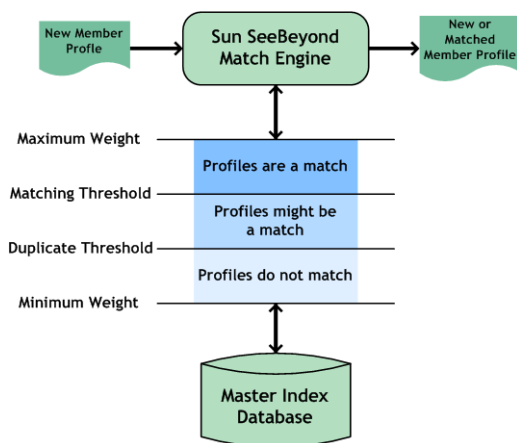


FIGURE 5 Weight Thresholds

Specifying the Weight Thresholds

There are many techniques for determining the initial settings for the match and duplicate thresholds. This section discusses two methods. You can also use the Data Profiler and Bulk Matcher to determine these thresholds. For more information, see [Analyzing and Cleansing Data for Sun Master Index](#) and [Loading the Initial Data Set for a Sun Master Index](#).

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 6](#). 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 6](#). 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.

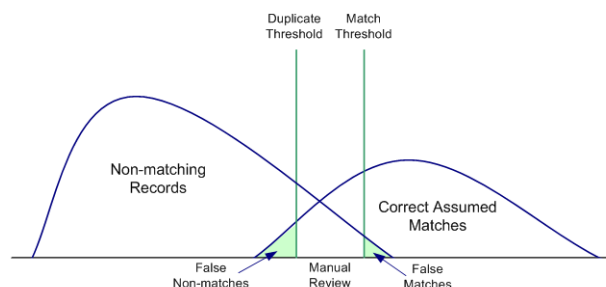


FIGURE 6 Weight Distribution Chart

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 62](#) or [“Weight Ranges Using Probabilities” on page 63](#), 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 17](#), 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 described earlier, 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.

