

# **Oracle Commerce**

**Content Acquisition System Relationship Discovery Guide**

**Version 11.2 • October 2015**





# Contents

<b>Preface.....</b>	<b>7</b>
About this guide.....	7
Who should use this guide.....	7
Conventions used in this guide.....	7
Contacting Oracle Support.....	7
 <b>Chapter 1: Introduction to Term Discovery.....</b>	 <b>9</b>
Overview of Relationship Discovery.....	9
Overview of Term Discovery.....	9
Overview of Cluster Discovery.....	11
 <b>Chapter 2: Configuration Guidelines for Term Extraction.....</b>	 <b>13</b>
Adding a Term Extractor manipulator.....	13
Term extraction workflow.....	14
Minimal term extraction configuration.....	14
Source input text property.....	15
Record specifier property name.....	15
Noun phrase grouping.....	16
Terms output property.....	17
All-terms destination property.....	17
Language.....	18
Supported languages.....	18
Configuration for the exclude list.....	19
Configuration for the main term extraction module.....	20
Update mode.....	20
Maximum number of input records.....	21
Configuration for candidate term identification.....	21
Input term property.....	22
Language specification of input records.....	22
Configuration for corpus-level filtering.....	23
Minimum and maximum occurrences in records.....	23
Minimum and maximum coverage settings.....	24
Threshold for the global informativeness of terms.....	24
Using regular expressions.....	24
Enabling debugging information for corpus-level filtering.....	25
Configuration for record-level filtering.....	25
Specifying a scoring threshold.....	26
Limiting the number of terms per record.....	26
Best practices for term filtering.....	27
Format of the source data.....	28
 <b>Chapter 3: Configuration Guidelines for Clustering.....</b>	 <b>31</b>
Configuration for clusters.....	31
Clustering parameter descriptions.....	32
Tuning strategy for clusters.....	33
 <b>Chapter 4: Building the Front End of the Term Discovery Application.....</b>	 <b>35</b>
Files to be changed.....	35
Adding global constants.....	35
Setting refinements in the controller file.....	36
Displaying refinements.....	37
Displaying clusters.....	37
Cluster properties.....	37
JSP code for displaying clusters.....	37
Clustering overlap properties.....	40
Displaying records and dimension refinements.....	40

<b>Chapter 5: Term Discovery Advanced Topics.....</b>	<b>41</b>
Term filtering with pre-tagged records.....	41
Filtering only pre-existing terms.....	41
Filtering both sets of terms uniformly.....	42
Filtering only the new terms.....	42
Tuning aids for the filtering parameters.....	43
Using STATEFUL mode for tuning.....	43
Using corpus-filtering logging statistics.....	44
<b>Appendix A: Term Discovery Sample Files.....</b>	<b>45</b>
Modified nav_controls.jsp file.....	45
New nav_clusters.jsp file.....	49

---

# Copyright and disclaimer

Copyright © 2003, 2015, Oracle and/or its affiliates. All rights reserved.

This software and related documentation are provided under a license agreement containing restrictions on use and disclosure and are protected by intellectual property laws. Except as expressly permitted in your license agreement or allowed by law, you may not use, copy, reproduce, translate, broadcast, modify, license, transmit, distribute, exhibit, perform, publish, or display any part, in any form, or by any means. Reverse engineering, disassembly, or decompilation of this software, unless required by law for interoperability, is prohibited.

The information contained herein is subject to change without notice and is not warranted to be error-free. If you find any errors, please report them to us in writing.

If this is software or related documentation that is delivered to the U.S. Government or anyone licensing it on behalf of the U.S. Government, the following notice is applicable:

U.S. GOVERNMENT END USERS: Oracle programs, including any operating system, integrated software, any programs installed on the hardware, and/or documentation, delivered to U.S. Government end users are "commercial computer software" pursuant to the applicable Federal Acquisition Regulation and agency-specific supplemental regulations. As such, use, duplication, disclosure, modification, and adaptation of the programs, including any operating system, integrated software, any programs installed on the hardware, and/or documentation, shall be subject to license terms and license restrictions applicable to the programs. No other rights are granted to the U.S. Government.

This software or hardware is developed for general use in a variety of information management applications. It is not developed or intended for use in any inherently dangerous applications, including applications that may create a risk of personal injury. If you use this software or hardware in dangerous applications, then you shall be responsible to take all appropriate fail-safe, backup, redundancy, and other measures to ensure its safe use. Oracle Corporation and its affiliates disclaim any liability for any damages caused by use of this software or hardware in dangerous applications.

Oracle and Java are registered trademarks of Oracle and/or its affiliates. Other names may be trademarks of their respective owners.

Intel and Intel Xeon are trademarks or registered trademarks of Intel Corporation. All SPARC trademarks are used under license and are trademarks or registered trademarks of SPARC International, Inc. AMD, Opteron, the AMD logo, and the AMD Opteron logo are trademarks or registered trademarks of Advanced Micro Devices. UNIX is a registered trademark of The Open Group.

This software or hardware and documentation may provide access to or information on content, products, and services from third parties. Oracle Corporation and its affiliates are not responsible for and expressly disclaim all warranties of any kind with respect to third-party content, products, and services unless otherwise set forth in an applicable agreement between you and Oracle. Oracle Corporation and its affiliates will not be responsible

---

for any loss, costs, or damages incurred due to your access to or use of third-party content, products, or services, except as set forth in an applicable agreement between you and Oracle.

For information about Oracle's commitment to accessibility, visit the Oracle Accessibility Program website at <http://www.oracle.com/pls/topic/lookup?ctx=acc&id=docacc>.

Oracle customers that have purchased support have access to electronic support through My Oracle Support. For information, visit <http://www.oracle.com/pls/topic/lookup?ctx=acc&id=info> or visit <http://www.oracle.com/pls/topic/lookup?ctx=acc&id=trs> if you are hearing impaired.

# Preface

Oracle Commerce Guided Search is the most effective way for your customers to dynamically explore your storefront and find relevant and desired items quickly. An industry-leading faceted search and Guided Navigation solution, Guided Search enables businesses to influence customers in each step of their search experience. At the core of Guided Search is the MDEX Engine™, a hybrid search-analytical database specifically designed for high-performance exploration and discovery. The Oracle Commerce Content Acquisition System provides a set of extensible mechanisms to bring both structured data and unstructured content into the MDEX Engine from a variety of source systems. The Oracle Commerce Assembler dynamically assembles content from any resource and seamlessly combines it into results that can be rendered for display.

Oracle Commerce Experience Manager enables non-technical users to create, manage, and deliver targeted, relevant content to customers. With Experience Manager, you can combine unlimited variations of virtual product and customer data into personalized assortments of relevant products, promotions, and other content and display it to buyers in response to any search or facet refinement. Out-of-the-box templates and experience cartridges are provided for the most common use cases; technical teams can also use a software developer's kit to create custom cartridges.

## About this guide

This guide describes the tasks involved in creating an Oracle Commerce Guided Search Relationship Discovery application.

## Who should use this guide

This guide is intended for developers responsible for creating a Relationship Discovery implementation.

## Conventions used in this guide

This guide uses the following typographical conventions:

Code examples, inline references to code elements, file names, and user input are set in `monospace` font. In the case of long lines of code, or when inline monospace text occurs at the end of a line, the following symbol is used to show that the content continues on to the next line: ~

When copying and pasting such examples, ensure that any occurrences of the symbol and the corresponding line break are deleted and any remaining space is closed up.

## Contacting Oracle Support

Oracle Support provides registered users with answers to implementation questions, product and solution help, and important news and updates about Guided Search software.

You can contact Oracle Support through the My Oracle Support site at <https://support.oracle.com>.





## Chapter 1

---

# Introduction to Term Discovery

This section provides overviews of the Guided Search Term Discovery and Guided Search Cluster Discovery features.

## Overview of Relationship Discovery

The Guided Search Relationship Discovery feature provides you with the ability to identify and extract key relationships in documents, including documents consisting of unstructured text.

This guide describes how to implement the Guided Search Term Discovery and Guided Search Cluster Discovery features, which are two major components of the Relationship Discovery solution. The third major component, Guided Search Entity Discovery, is not documented in this guide; for information on this feature, contact your Oracle Commerce Guided Search representative.



**Note:** The implementation of Term Discovery and Cluster Discovery requires that you must have purchased a license for Relationship Discovery. The license package includes documentation that describes how to enable the feature. Contact your Oracle Commerce Guided Search representative if you are not certain whether you have a Relationship Discovery license.

The Term Discovery and Cluster Discovery components of Relationship Discovery have the capability to:

- Extract salient terms (noun phrases) from documents.
- Provide a scoring mechanism for the extracted terms, which determines whether a given term is retained. Note that term scores are for internal use by the Guided Search software and therefore not exposed.
- Generate relevant terms on queries.
- Generate clusters of relevant terms.

## Overview of Term Discovery

Term Discovery is the feature that extracts salient terms from source documents.

Term Discovery can be thought of as a two-part process:

1. Extracting terms from source documents (unstructured or structured) and scoring them according to their relevancy. The scored terms are mapped to a Guided Search dimension, called a *Term Discovery dimension* in this guide.
2. Presenting terms relevant to the current navigation state.

## Extracting terms from documents

*Term extraction* is the process of tagging a Guided Search record with a list of its relevant terms. A *term* represents a concept mentioned in the record's source document, and is typically a noun phrase. The noun phrase consists of one or more nouns, potentially with adjoining words. A relevant term is a term that bears information for a document relative to the rest of the corpus.

During the term extraction process, term variants found in documents are stemmed for comparison and aggregation, but the final representation of the term uses the dominant form (most frequent variant). Using the dominant form allows the recovery of the preferred representation (singular/plural case, capitalization) of proper nouns and brand names.

Term extraction is performed by CAS using a CAS term extractor manipulator. Configuration information on term extraction is given in the *Configuration Guidelines for Term Extraction* chapter.

## Maximum size of extracted terms

A noun phrase consists of a noun (or a sequence of nouns) with any associated modifiers. The modifiers are limited to adjectives, adjective phrases, or nouns used as adjectives.

Programmatically, each word in a noun phrase is called a token. An extracted noun phrase can have a maximum of 5 tokens; each token is limited to 200 characters. Therefore, a valid noun phrase has a maximum size of 1,000 characters.

The maximum size of a sentence in a document is 1,000 tokens (words). If the term extractor cannot determine the sentence boundaries of text in the document, it splits the text into blocks of 1,000 tokens and then performs text extraction on the blocks. In addition, the following WARN message will be entered in the term extraction log:

```
Sentence boundaries could not be found for text beginning
with tokens t1 t2 t3 t4 t5
```

where t1 through t5 are the first 5 tokens of the problematic text.

If you are using OLT noun grouping, the rules in the list below are applied only when the noun group size is more than 1000 characters. If the length of the noun group is less than or equal to 1000 characters, the noun group is taken as is without any modification.

The term extractor treats invalid noun phrases as follows:

- If a noun phrase has more than 5 tokens, only the last 5 tokens are retained. For example, with a 7-token noun phrase, the first 2 tokens are completely ignored by the term extractor and the last 5 tokens are retained.
- Tokens that are over 200 characters are ignored.
- If a noun phrase includes an overly-long token, that token is ignored, and the precedent and antecedent tokens are treated as separate noun phrases. For example, assume a 5-token noun phrase. Token 3 is an overly-long token and the others are valid. In this case, Token 3 will be ignored and the term extractor will return 2 noun phrases: the first noun phrase will consist of Tokens 1 and 2, and the second will consist of Tokens 4 and 5.

## Presenting relevant terms

*Relevant terms* are the most frequent terms available in the Term Discovery dimension. These terms are returned from the documents in the current result set. All of the terms in the set belong to the same Term Discovery dimension.

The Term Discovery dimension must have these two attributes:

- It must be a flat dimension (that is, a dimension that does not contain hierarchies).

- It must not be a hidden dimension. (Configuring it as a hidden dimension will disable the Cluster Discovery feature.)

Relevant terms are returned by the MDEX Engine as dimension value refinements. Programmatically, relevant terms are `DimVal` objects. Therefore, application developers can use the same Presentation API methods on relevant terms that can be used on normal dimension value refinements. For example, they can be returned sorted using any ranking behavior supported for dimension value refinements.

## Overview of Cluster Discovery

A cluster is a collection of relevant terms, providing a grouping of Guided Search records that share these common terms.

All of the terms, which are dimension values, must come from the same dimension and must be a Term Discovery dimension. Clusters can be generated only if the Term Discovery dimension is available for refinements. This dimension cannot be hidden, and it must also be available from the navigation states for which you want clusters. Therefore, your application must have this dimension globally available, instead of having it available only when triggered by precedence rules.

The following features apply to the clusters:

- The MDEX Engine performs dynamic clustering. That is, when a user navigates the clustering tree, it is re-clustered at any selection, allowing users to zoom into their data to practically any level.
- There is no limit in the number of records that can be clustered.
- Each cluster is represented by a list of terms, which provide to the user what is known as *information scent*: the user is instantly aware of what each cluster contains. The user can quickly understand the implied content of the clustered records.
- All clusters are designed to maximize two metrics: coherence -- each cluster has only closely related records -- and distinctiveness -- two different clusters will have different records.
- Each cluster has high recall. A match partial technique is typically used on cluster selection, maximizing the number of semantically related records that are returned.

For more information, see the *Configuration Guidelines for Clustering* chapter.



## Chapter 2

---

# Configuration Guidelines for Term Extraction

This topic describes some guidelines for configuring term extraction.

## Adding a Term Extractor manipulator

You add a **Term Extractor** manipulator to a data source on the **Edit** page of CAS Console. The manipulator extracts terms from a given property on each record.

*Term extraction* is the process of tagging a record with a list of its relevant terms. A term represents a concept mentioned in the record's source document, and is typically a noun phrase. The noun phrase consists of one or more nouns, potentially with adjoining words. In CAS, the Term Extractor manipulator uses the following Oracle Language Technology (OLT) features to analyze a specific record property:

- Sentence recognition  
Identifies breaks between sentences in text.
- Tokenization  
Breaks a stream of text up into words, phrases, symbols, or other meaningful elements.
- Dynamic stemming  
Determines the base (uninflected) form of a word. The process is based on dictionary entries and language-specific rules.
- Parts of speech tagging  
Tags a word in a text as belonging to a particular part of speech, based on its definition and its context.
- Noun phrase grouping  
Identifies noun phrases in text using grammatical phrase grouping. You can use either OLT or custom technology.

The terms are extracted into a property that you specify in the record. To add a Term Extractor manipulator to a data source:

1. On the **Data Sources** page, click a data source name to access its acquisition steps.
2. Click **Add Manipulator...** and select **Term Extractor**.
3. Click **Add**.  
The **Manipulator Settings** page displays.

4. In **Name**, specify a unique name for the manipulator to distinguish it from other manipulators in the data source.

You can specify a name with alphanumeric characters, underscores, dashes, and periods. All other characters are invalid for a name. If you change this value, you must run a full crawl.

5. Specify the following required configuration parameters:

- TEXT\_PROP\_NAME
- RECORD\_SPEC\_PROP\_NAME
- Use OLT for NounGrouping
- OUTPUT\_PROP\_NAME
- ALL\_TERMS\_OUTPUT\_PROP\_NAME
- LANG
- State Directory path

If you modify any of these parameters except for LANG, you must run a full crawl. For information on these parameters, see *Minimum term extraction configuration*.

6. You can specify an optional pass-through parameter.
7. Click **Save**.

The manipulator displays on the **Acquisition Steps** list.

## Term extraction workflow

This topics gives an overview of the workflow of the term extractor.

Term extraction consists of three steps, each of which is optional:

1. Candidate Term Identification — Identify all terms that are candidates for a given document and then extract those terms. This step is omitted if terms are being extracted from pre-tagged records.
2. Corpus-level Filtering — Globally filter the extracted terms to determine a controlled vocabulary for the corpus. This involves identifying terms that should be removed corpus-wide.
3. Record-level Filtering — Determine, for each record, what are the most relevant terms for it from the controlled vocabulary. This involves identifying terms that should be removed from an individual record (independent of terms that should be removed from the entire corpus, but possibly using corpus-level information) and subsequently removing these terms from the tags on the record.

Note that steps 2 and 3 remove the terms that are judged by their score or by their presence on the exclude list to be of low information value. The tagged terms on each record are a result of step 3.

This section provides configuration requirements for the term extraction modules. You supply the configuration parameters as pass-through name/value pairs to the CAS manipulator.

## Minimal term extraction configuration

The minimal term extraction configuration consists of the following required pass-through parameters.

These parameters are listed in the following table, with details in later sections.

Parameter	Configuration Value
TEXT_PROP_NAME	Source property to use as the source text for term extraction. No default.
RECORD_SPEC_PROP_NAME	The source property that is mapped to the Guided Search record specifier property. No default.
Use OLT for NounGrouping	Property determines if term extraction uses OLT (true) or legacy code (false) for noun phrase grouping.
OUTPUT_PROP_NAME	The source property to use as the destination for tagged terms. No default.
ALL_TERMS_OUTPUT_PROP_NAME	Source property to use as the destination for all terms on a record. No default.
LANG	Specifies the language ID to use on a global basis. No default.
State Directory path	Directory path where term extraction state files should be stored. The location should be Guide Search application specific. For example: <ENDECA_APPS_DIR>/Discover/state/TE..

This configuration runs as follows:

- A baseline update (STATELESS mode) is performed.
- No corpus- or record-level filtering will be performed. As a result, all terms in the corpus are extracted and all of them are tagged onto records.
- The expected language of the documents are specified in the LANG pass-through.

This configuration is the most permissive for term extraction. Although most sites prefer to perform some level of corpus and record filtering, this minimal configuration might be useful with small data sets that have a closely-related set of noun phrases in their documents.



**Note:** You must run a full crawl if any one of the minimum configuration parameters except LANG are modified.

## Source input text property

The TEXT\_PROP\_NAME pass-through specifies which pre-existing property on the input data record will be used as the source for term extraction.

The TEXT\_PROP\_NAME pass-through is mandatory. Note that the value for this pass-through specifies a source property, not a Guided Search property. This property will then be mapped to a dimension.

If you want to extract terms from multiple properties, a pipeline component must combine the text from the multiple properties into one text property. The name of that text property is then used for the TEXT\_PROP\_NAME pass-through.

## Record specifier property name

The RECORD\_SPEC\_PROP\_NAME pass-through specifies a source record property that is mapped to the Endeca record specifier property in the implementation.

The RECORD\_SPEC\_PROP\_NAME pass-through is mandatory. Get the RECORD\_SPEC\_PROP\_NAME pass-through from the `<ENDECA_APP>\config\cas\data-recordstore-config.xml` file. This file is used to set the definition for the Record Store and it contains the RECORD\_SPEC\_PROP\_NAME property.

To find out the name of the record specifier source property:

1. Go to `<ENDECA_APP>\config\cas\data-recordstore-config.xml` file
2. Open the file in an editor.

For example, in the Discover reference application:

```
<recordStoreConfiguration xmlns="http://recordstore.itl.endeca.com/">
  <idPropertyName>common.id</idPropertyName>
</recordStoreConfiguration>
```

The `common.id` property is the property name that you must include in the RECORD\_SPEC\_PROP\_NAME pass-through of the CAS manipulator.

3. Use this name for the RECORD\_SPEC\_PROP\_NAME pass-through.

## Noun phrase grouping

You can use either custom noun phrase grouping or OLT for noun phrase grouping through the CAS manipulator configuration.

In both options, OLT is used for tokenization, part-of-speech tagging, dynamic stemming, and sentence recognition.

For the **Use OLT for NounGrouping** setting:

- A **false** value means custom logic is used for noun phrase grouping. This model helps you get similar behavior as in Forge-based term extraction.



**Note:** Only English and French languages are supported by this model.

- A **true** value means that the OLT grammatical noun phrase grouping is used in this module. This model supports over fifty languages including English and French. See [Supported languages](#) on page 18 for more information on the supported languages.

## Output comparison

This section compares the output of a false value (Forge-based term extraction) and a true value (OLT-based term extraction in CAS). In OLT-based term extraction, noun groups are formed with better meaning and context. The following example uses the same text given in both Forge-based term extraction and OLT-based term extraction.

Here is the text:

Sachin Ramesh Tendulkar is a former Indian cricketer widely acknowledged as the greatest batsmen of all time, popularly holding the title "God of Cricket" among his fans. He is the only player to have scored one hundred international centuries.

The following table shows the noun groups that were extracted:

Forge-based term extraction	OLT-based term extraction using CAS
Sachin Ramesh Tendulkar	Sachin Ramesh Tendulkar
cricketer	a former Indian cricketer



Forge-based term extraction	OLT-based term extraction using CAS
greatest batsmen	the greatest batsmen of all time
title	the title
God	God of Cricket
fans	his fans

## Terms output property

The `OUTPUT_PROP_NAME` pass-through specifies the property that holds the tagged terms (if any) for each Endeca record.

The `OUTPUT_PROP_NAME` pass-through is mandatory. If the property does not already exist in the implementation, then the term extractor will create it. After term extraction, this property is typically mapped to a dimension via the property mapper.

## All-terms destination property

The `ALL_TERMS_OUTPUT_PROP_NAME` pass-through specifies the property which gets all terms on a record that pass corpus-level filtering.

The `ALL_TERMS_OUTPUT_PROP_NAME` pass-through is mandatory. The property gets all terms on a record that pass corpus-level filtering regardless of whether they pass record-level filtering.

If the property does not already exist in the implementation, the term extractor will create it. When mapped to a Guided Search property, record searches can be performed on this all-terms property (assuming it is configured to be searchable).

The all-terms property is used for search purposes. For each record, the term extractor finds all of the terms in the corpus-wide vocabulary that occur in that document (regardless of their relevance for that document) and puts them in the all-terms property. By using this property for searches, stemming of the terms can be performed.

For example, if the terms "search engine" and "search engines" appear in the corpus, they will be normalized to the dominant form (e.g., as "search engine"). But if you want to find all records that contain either variant, you cannot use Phrase search against the body text, because Phrase search does not locate stemmed variants. Instead, Term Discovery ensures that both the dimension value name and the term stored in the all-terms property is the dominant form.

The terms in the all-terms property are separated by using **sep** as the delimiter. The term extractor makes the separator by doing `(sep)+` until the separator is not a substring of any term in the corpus. Therefore, the separator may be **sep**, **sepsep**, **sepsepsep**, and so on. For example, an article on Aldera, Spain might produce the following all-terms property (named `P_AllTerms`) on the record (in the example, **sepsep** is the separator for the property):

```
P_AllTerms: district sepsep Spain sepsep south coast sepsep
            coast sepsep town sepsep Aldera sepsep province sepsep
            Romans sepsep decline sepsep station sepsep hills sepsep
            Heracles sepsep temple sepsep colonies sepsep Tiberius
```

Because of the widespread use of this separator, you should add it to the stop word list. (Note that this is the application's stop word list, not the term exclude list.) Before doing so, first determine which form is the separator. Run the corpus at least once to find what the separator is and then set that separator as a stop word. For example, if "sep" is a valid term in the corpus, then it is likely that **sepsep** will be the separator. Thus, you would

add "sepsep" (but not "sep") to the stop word list. Then, periodically monitor the corpus to make sure the separator has not changed.

## Language

The LANG pass-through parameter sets the language ID of input records on a global basis

The language ID is not case sensitive. For example, you can specify EN or en for English (American).

## Supported languages

The LANG pass-through parameter takes one of the following language codes as input when **Use OLT for noun grouping** is true.

Language	Language Code
Arabic	ar
Basque	eu
Belarusian	be
Bosnian	bs
Bulgarian	bg
Catalan	ca
Chinese (Simplified)	zh-CN
Chinese (Traditional)	zh-TW
Croatian	hr
Czech	cs
Danish	da
Dutch	nl
English (American)	en
English (United Kingdom)	en-GB
Estonian	et
Farsi/Persian	fa
Finnish	fi
French (European)	fr
French (Canadian)	fr-CA
Galician	gl
German	de
Greek	el
Hebrew	he

Language	Language Code
Hungarian	hu
Indonesian	id
Italian	it
Japanese	jp
Korean	ko
Latvian	lv
Lithuanian	lt
Macedonian	mk
Malay	ms
Norwegian (Bokmal)	nb
Norwegian (Nyorsk)	nn
Polish	pl
Portuguese (European)	pt
Portuguese (Brazil)	pt-BR
Romanian	ro
Russian	ru
Serbian (Cyrillic)	sr
Serbian (Latin)	sr-Latn
Slovak	sk
Slovenian	sl
Spanish	es
Swedish	sv
Thai	th
Turkish	tr
Ukranian	uk
Valencian	ca-ES-valencia
Vietnamese	vi

## Configuration for the exclude list

The exclude list configuration contains a set of terms that are removed from the final list of extracted terms.

Excludes are compared against the canonical and all raw forms of a term; if it matches any, the term is excluded. This is equivalent to canonicalizing the exclude term.

The exclude list configuration can be passed to the CAS term extraction manipulator by creating a new Record Store of a supported type, for example de-limited or JDBC. You must load the data into the Record Store and add the following pass through information in the manipulator configuration.

Parameter	Configuration Value
Exclude List Record Store instance name	Record Store instance name which contains exclude terms.
Exclude term property name	Property name of the record which contains exclude term in Exclude List Record Store.

The format rules for the excluded terms list are as follows:

- The exclude list is not case sensitive.
- Apostrophes can be included as necessary. For example, enter `i'm` if you want that term to be excluded.

The list is processed after all terms have been extracted from the records.

In this brief example of a delimited exclude list file, `EXCLUDE` and `RecordSpec` are the headers -- multiple headers are allowed. For the `Exclude term property name` property in the CAS term extraction manipulator, you must pass `EXCLUDE`.

```
EXCLUDE,RecordSpec
- 12.1 megapixel,1
- 12 MP,2
The 12.1 megapixel sensor,4
```

## Configuration for the main term extraction module

In addition to the parameters in the minimum configuration, this module requires the parameters listed in the following table.

PASS_THROUGH Element	Configuration Value
UPDATE_MODE	The type of data update to perform: <code>STATELESS</code> (the default), <code>STATEFUL</code> , or <code>PARTIAL</code> . Optional.
MAX_INPUT_RECORDS	Integer that sets the maximum number of records to be processed. Optional. Defaults to all records processed.

### Update mode

The `UPDATE_MODE` pass-through specifies which type of data update is being performed by the pipeline.

The `UPDATE_MODE` pass-through is optional. The three values for this pass-through are `STATELESS`, `STATEFUL`, and `PARTIAL`. Note that if this pass-through is omitted, the term extractor performs a `STATELESS` update.

#### STATELESS mode

`STATELESS` is analogous to a baseline update and performs the following actions:

1. Extracts terms from all records.
2. Performs corpus-level and record-level filtering on all records.
3. Emits all records.

4. Does not create state files.

### STATEFUL mode

STATEFUL performs the following actions:

1. Extracts terms from new records only.
2. Performs corpus-level and record-level filtering on all records -- both previously processed records and new records.
3. Emits all records -- both previous and new records.
4. Creates term state files.

### PARTIAL mode

PARTIAL is analogous to a partial update and performs the following actions:

1. Extracts terms from new records only.
2. Performs record-level filtering on new records only. Corpus-level filtering is not performed at all, so the previous corpus information is left as-is.
3. Emits new records only.
4. Does not create state files. Previous term state files are left as-is.

### Notes on update modes

Keep the following notes in mind when using the update modes:

- STATELESS mode is recommended for sites that perform baseline updates exclusively.
- STATEFUL mode is recommended for sites that implement partial updates. That is, the baseline update pipeline will use STATEFUL mode while the partial update pipeline will use PARTIAL mode.
- The STATELESS and STATEFUL modes do not need pre-existing state files in order to run.
- The PARTIAL mode does require the existence of the term state files. Therefore, a STATEFUL update must be performed before a PARTIAL update.

## Maximum number of input records

The MAX\_INPUT\_RECORDS pass-through specifies the maximum number of input records to be processed. With the MAX\_INPUT\_RECORDS pass-through, only the specified number of records are read in and emitted.

The MAX\_INPUT\_RECORDS pass-through is optional. This pass-through is intended for development purposes and should be omitted in a production environment so that all records are processed.

## Configuration for candidate term identification

A set of pass-throughs is available to configure which terms are candidates for term extraction.

The following configuration parameters are used in addition to the parameters in the minimum configuration when performing candidate term identification for terms:

PASS_THROUGH Element	Configuration Value
INPUT_TERM_PROP_NAME	Source property to use as the source text for pre-tagged records. Optional. No default.

PASS_THROUGH Element	Configuration Value
LANG_PROP_NAME	Source property containing the language ID. Optional. No default.
LANG	Specifies the language ID to use on a global basis. No default.

## Input term property

The INPUT\_TERM\_PROP\_NAME pass-through is used for a corpus that contains pre-tagged records.

Pre-tagged records are source records with pre-existing terms that were generated by non-Guided Search software. In this scenario, you do not want to extract new terms from the documents but do want to perform corpus-level and/or record-level filtering on the pre-tagged terms.

This pass-through can be used in conjunction with the TEXT\_PROP\_NAME pass-through to combine the pre-existing terms with the new extracted terms.

## Language specification of input records

Two pass-through parameters set the language ID of input records on a global and per-record basis.

You can use the LANG and LANG\_PROP\_NAME pass-through parameters to specify the global language ID and the per-record language ID of the input records. The language ID is not case sensitive for both pass-through parameters. For example, you can specify `EN` or `en` for English (American).

Note that the LANG\_PROP\_NAME value takes precedence, and if not present, the value of LANG is used as the language of the record.

Both can be specified in the CAS manipulator.

### LANG pass-through

The LANG pass-through specifies the language ID to use on a global basis.

### LANG\_PROP\_NAME pass-through

The LANG\_PROP\_NAME pass-through specifies the name of the record property that contains the language ID for that record. If you do not specify this pass-through, the language ID for each record will default to the value of the LANG pass-through. For example, if the value for LANG is `en-GB`, then the term extractor assumes that all the records are in English - UK.

If you do specify the LANG\_PROP\_NAME pass-through, the term extractor will evaluate each record as follows:

- If the value of the LANG\_PROP\_NAME property matches the LANG setting, then terms are extracted from the record in that language.
- If the value of the LANG\_PROP\_NAME property does not match the LANG setting, then terms are not extracted from the record. That is, the record is ignored for purposes of term extraction, but the record is otherwise processed by CAS. For example, if the value of LANG is `fr` (French - European) and the value of the LANG\_PROP\_NAME property is `en-GB` (English - UK) for a given record, the terms extractor will ignore that record.
- If the value of the LANG\_PROP\_NAME property is null or the record does not contain the LANG\_PROP\_NAME property, the term extractor will assume that the language ID of the record is the same as the LANG setting and therefore will attempt to extract terms from the record in that language.

If you have documents in multiple languages, the LANG\_PROP\_NAME pass-through is useful to ensure that only records in the desired language (the LANG setting) are processed by the term extractor.

## Configuration for corpus-level filtering

A set of pass-through parameters is available to configure how corpus-level filtering is done.

In addition to the parameters in the minimum configuration, the following configuration parameters are used for performing corpus-level filtering.

PASS_THROUGH Element	Configuration Value
CORPUS_MIN_RECS	Integer specifying the minimum number of records in which a term must appear in order to be considered. Default is 0; minimum recommended value is 2.
CORPUS_MAX_RECS	Integer specifying the maximum number of records in which a term must appear in order to be considered. Default is unlimited.
CORPUS_MIN_COVERAGE	Double specifying the minimum fraction of the corpus that must contain the term. Default is Double.NEGATIVE_INFINITY; useful range is 0–1; recommended value is 0.00005.
CORPUS_MAX_COVERAGE	Double specifying the maximum fraction of the corpus that must contain the term. Default is Double.POSITIVE_INFINITY; useful range is 0–1; recommended value is 0.2.
CORPUS_MIN_INFO_GAIN	Minimum global informativeness score a term must have to be considered. Default is Double.NEGATIVE_INFINITY.
CORPUS_MAX_INFO_GAIN	Maximum global informativeness score a term must have to be considered. Default is Double.POSITIVE_INFINITY.
CORPUS_REGEX_KEEP	String that is a regular expression for terms that should be kept. No default.
CORPUS_REGEX_SKIP	String that is a regular expression for terms that should be discarded. No default.
CORPUS_DEBUG	If set to <code>true</code> , detailed information about corpus-level scoring is written to the logs. Default is <code>false</code> .

### Minimum and maximum occurrences in records

Two pass-throughs set the minimum and maximum number of records in a term can appear before it is discarded.

The `CORPUS_MIN_RECS` parameter determines the minimum number of records in which a term can appear before it can be considered. If a term appears in fewer records than the `CORPUS_MIN_RECS` setting, then the term is discarded.

Setting a value less than 2 is not recommended, as it is useless for clustering. In addition, a value of 2 or higher means that singletons are eliminated.

Singletons are terms that appear in only one record. Singletons typically occur very frequently in a corpus, but cannot be used for clustering, which uses cross-record statistics. Therefore, eliminating singletons reduces memory use and computation time. Minimum occurrences should be set to at least 2 to remove singletons; higher values will require a term to appear on more records. Note that large document sets have more term redundancy and can have this parameter set to higher values, such as 3 or higher.

The `CORPUS_MAX_RECS` parameter sets the maximum number of records in which a term can appear. If it appears in more records than the set value, then the term is discarded.

You can use both pass-throughs to create a window. For example, assume you have set these values:

```
CORPUS_MIN_RECS=5
CORPUS_MAX_RECS=20
```

The result would be that a term would be retained only if it appeared in at least 5 records but no more than 20 records.

## Minimum and maximum coverage settings

Two pass-throughs set the minimum and maximum percentage of records that must contain a term.

The coverage pass-throughs correspond to the fraction of the records in the corpus which contain at least one occurrence of the given term:

- `CORPUS_MIN_COVERAGE` sets the minimum fraction of the corpus that must contain the term.
- `CORPUS_MAX_COVERAGE` sets the maximum fraction of the corpus that must contain the term.

For example, if you want to keep only those terms that appear in between 5% and 25% of the corpus, you would use these settings:

```
CORPUS_MIN_COVERAGE=0.05
CORPUS_MAX_COVERAGE=0.25
```

The useful range for both pass-throughs is 0-1, in which 1 is 100% of the corpus.

## Threshold for the global informativeness of terms

Two pass-throughs set the scoring threshold for the global informativeness of a term.

The `CORPUS_MIN_INFO_GAIN` parameter sets the minimum scoring threshold when measuring the global informativeness (`info_gain`) of a term. The useful range for `CORPUS_MIN_INFO_GAIN` is minus infinity to plus infinity, although most terms tend to fall in the -5.0...+5.0 range. The terms with negative information gain are likelier to contribute more noise than signal. Eliminating these globally uninformative terms saves considerable memory and query compute time. The minimum `info_gain` setting can be increased to require higher global informativeness—or, conversely, decreased to allow more terms to be retained. Setting `CORPUS_MIN_INFO_GAIN` to 0 is usually adequate. Setting this parameter to values higher than 1.0-2.0 can dramatically reduce the number of terms.

The `CORPUS_MAX_INFO_GAIN` parameter sets the maximum scoring threshold for the global informativeness of a term. If the scoring for a term exceeds this threshold, the term is discarded.

## Using regular expressions

Two pass-throughs allow you to use regular expressions to retain or discard terms.

The term extractor supports two pass-throughs that are used for matching character sequences against patterns specified by regular expressions:

- `CORPUS_REGEX_KEEP` specifies a regular expression that a term must match in order to be retained.
- `CORPUS_REGEX_SKIP` specifies a regular expression that a term must not match in order to be retained.

You can use either or both pass-throughs. If both are used, then a term must pass both tests in order to be retained (that is, if the term satisfies one pass-through but not the other, the term is discarded).



The term extractor implements Oracle's `java.util.regex` package to parse and match the pattern of the regular expression. Therefore, the supported regular-expression constructs are the same as those in the documentation page for the `java.util.regex.Pattern` class at this URL:

<http://docs.oracle.com/javase/8/docs/api/java/util/regex/Pattern.html>

This means that among the valid constructs you can use are:

- Escaped characters, such as `\t` for the tab character.
- Character classes (simple, negation, range, intersection, subtraction). For example, `[^abc]` means match any character except a, b, or c, while `[a-zA-Z]` means match any upper- or lower-case letter.
- Predefined character classes, such as `\d` for a digit or `\s` for a whitespace character.
- POSIX character classes (US-ASCII only), such as `\p{Alpha}` for an alphabetic character, `\p{Alnum}` for an alphanumeric character, and `\p{Punct}` for punctuation.
- Boundary matchers, such as `^` for the beginning of a line, `$` for the end of a line, and `\b` for a word boundary.
- Logical operators, such as `X|Y` for either X or Y.

For a full list of valid constructs, see the `Pattern` class documentation page at the URL listed above.

The following is an example of a useful regular expression that uses the POSIX `\p{Alnum}` construct:

```
^\p{Alnum}[\p{Alnum}\.\- ' ]+$
```

When used with the `CORPUS_REGEX_KEEP` pass-through, this regular expression will retain only terms that have at least two characters, starts with an alphanumeric character, and contains only alphanumeric, period, dash, apostrophe, and space characters. (The apostrophe is to retain terms such as `O'Malley`).

The following is another example that is fairly restrictive:

```
[^A-Za-z0-9\-\.\ ]
```

When used with the `CORPUS_REGEX_SKIP` pass-through, this regular expression will retain only terms that consist of alphanumeric, dash, period, and space characters.

## Enabling debugging information for corpus-level filtering

The `CORPUS_DEBUG` pass-through enables the logging of debugging information.

The `CORPUS_DEBUG` pass-through enables the term extractor to write detailed information about the corpus-level filtering scores it assigns to terms. Temporarily setting this pass-through to `true` will help you to tune corpus-level filtering.

## Configuration for record-level filtering

Two pass-through parameters are available to configure record-level filtering for term extraction.

In addition to the parameters in the minimum configuration, the following configuration parameters are used for performing record-level filtering. These pass-through parameters are optional. Note that if you chose OLT noun phase grouping in the minimum configuration, you cannot use record-level filtering.

PASS_THROUGH Element	Configuration Value
RECORD_FRACT_OF_MEDIAN	Double that sets the minimum scoring fraction of the median to use. 1 . 1 is the recommended value; 0 . 0 is the default.
RECORD_NTERMS	Sets a limit on the maximum number of terms that are tagged on a record. Default is to have no limit.

## Specifying a scoring threshold

The `RECORD_FRACT_OF_MEDIAN` sets a scoring threshold for record-level filtering.

CAS uses a scoring method in which terms that are more frequent in this document than across the corpus are considered to be more relevant and thus are retained. Filtering occurs on the document-by-document basis; in other words, each term is considered for inclusion or exclusion separately for each document in which it occurs.

The distribution of scores for terms on a single record typically has very few high-scoring terms, followed by a long, gently-sloped plateau of marginally informative terms, with a sudden drop-off of few uninformative terms.

The `RECORD_FRACT_OF_MEDIAN` value lets you set a scoring threshold for the plateau; only terms that score above this threshold are kept. `RECORD_FRACT_OF_MEDIAN` should be set to a value that expresses the threshold as a fraction of the median score for terms on the document.

The recommended threshold is `1 . 1` (i.e., 10% higher than the median), which will keep only the highly-informative terms. Higher values will tend to increase precision (the terms that are kept are more likely to be relevant) but decrease recall (more likely to lose relevant terms). The default value of this threshold is `0 . 0`, which allows all terms through.

## Limiting the number of terms per record

The `RECORD_NTERMS` pass-through sets a limit on the maximum number of terms that are tagged on a record.

You can use the `RECORD_NTERMS` pass-through to implement one of two strategies to limit the number of terms that are tagged on records:

- Set an absolute upper limit.
- Establish a cut-off window.

You cannot mix both strategies. In both strategies, CAS determines which terms have the highest relevance for that record. Note that this pass-through is recommended mainly for collections that have large documents.

### Setting a hard limit

To set an absolute upper limit, use the `RECORD_NTERMS` pass-through with only one integer value. Use this version of the pass-through when you are certain about the number of terms you want per record and can therefore set a hard limit. In this example, `RECORD_NTERMS` is set to a value of 8:

Using this setting, CAS determines which are the eight most relevant terms for this record and tag the record with them.

### Establishing a cut-off window

To establish a cut-off window, use the `RECORD_NTERMS` pass-through with a range of two integers, which sets the lower and upper limits of a cut-off window. This windowing strategy establishes a window that will be scanned for an optimal cut-off. This cut-off is where term informativeness drops off most precipitously. Use this strategy when you want CAS to be sensitive to actual term informativeness rather than just using a hard limit.

You can think of the term range as providing a fuzzy neighborhood to be used instead of a hard limit. For example, instead of `RECORD_NTERMS` having a hard limit of 32, you can set it to a range of 24-36. This range establishes a window where a record can have a minimum of 24 terms and a maximum of 36 terms. CAS determines the optimal cut-off within that window for each record.

For example, assume that 40 terms were extracted from Record A and also from Record B:

- For Record A, the optimal cut-off for the terms might be after term 26 (because of a sharp drop-off in relevancy for terms 27-40). Therefore, Record A will have 26 terms tagged onto it.
- For Record B, the optimal cut-off for its set of terms might be after term 30. In this case, Record B will have 30 tagged terms.

When using the range version of this pass-through, keep the following in mind:

- The lowest recommended value for the lower limit is around 10. The reason is that the scores of the top terms scores usually differ noticeably, and the largest score drop-off is likely to be found at the setting for the lower limit. Thus, if the lower is less than 10, you should expect it to behave like the hard-limit version of RECORD\_NTERMS, which is misleading.
- The value for the upper limit should not be much larger than the value for the lower limit. If the difference is too much, the number of terms assigned to each particular record will be essentially random (within the cut-off window). The only way to have this number of terms relatively consistent is to use a lower- and upper-limit pair that are not too far from each other.

## Best practices for term filtering

This topic presents a best practices list for term extraction.

The best practices lists include the pass-through name and a recommended (best practices) value. If a recommended value is not given, then the default value is also the recommended one.

The two important things to keep in mind are:

- When tuning corpus-level filtering, the number of documents is the main consideration.
- When tuning record-level filtering, the size of the text property is the main consideration.

### Corpus-level filtering best practices

#### CORPUS\_MIN\_RECS

Recommendation: Values of less than 2 are not recommended in general, since they allow terms that are seen only once in the entire corpus. If clustering is used, this value **MUST** be set to at least 2. Note that this parameter works similarly to CORPUS\_MIN\_COVERAGE: terms that are seen less frequently than in CORPUS\_MIN\_RECS are discarded, as are terms that are seen in less than  $\text{CORPUS\_MIN\_COVERAGE} * (\text{number of documents in the corpus})$ .

#### CORPUS\_MAX\_RECS

Recommendation: As a general rule of thumb, this pass-through does not have to be used. If your number of records can change, for example, through partial updates, Oracle recommends that you not use CORPUS\_MAX\_RECS, because the statistics will change with the changed number of records. In this case, you may want to use the CORPUS\_MAX\_COVERAGE pass-through instead.

#### CORPUS\_MIN\_COVERAGE

Recommendation: The useful range is 0-1. A value of 0.00005 is a good compromise, because the term extractor will retain a term if it has been seen in at least one document out of 20,000.

This value will change with the nature of the data set. For example, a site with a data set with a lot of topical diversity (such as news) can reduce this value and allow terms with lower coverage (however, one out of any hundred thousand is probably the smallest reasonable value). If memory use is an issue, you should increase this value.

#### CORPUS\_MAX\_COVERAGE

Recommendation: The useful range is 0-1. A value of 0.2 (which is 20% of the documents) is a good compromise. If a term is seen in more than one out of five documents (i.e., 20%), it is probably too broad to be useful. If terms that are tagged onto documents seem too generic, this number should be turned down. As with `CORPUS_MIN_COVERAGE`, turning this number down, even slightly, should free memory.

### **CORPUS\_REGEX\_KEEP**

Recommendation: A useful regular expression for terms to keep is:

```
^\p{Alnum}[\p{Alnum}\.\- ' ]+$
```

This retains terms that have at least two characters, start with an alphanumeric character, and includes only alphanumerics, spaces, periods, dashes, and single quotes.



**Note:** Each term must both match `CORPUS_REGEX_KEEP` and not match `CORPUS_REGEX_SKIP` to be retained.

### **CORPUS\_REGEX\_SKIP**

Recommendation: Use this pass-through only if you are certain of the format of the terms you want to discard.

### **CORPUS\_MIN\_INFO\_GAIN and CORPUS\_MAX\_INFO\_GAIN**

Recommendation: Begin by setting `CORPUS_MIN_INFO_GAIN` to 0. Do not set `CORPUS_MAX_INFO_GAIN` initially. Tune the other term extraction pass-throughs. Then, run a data set (or a subset) with `CORPUS_DEBUG` set to `true`, which will print the list of terms that passed all the selection criteria. You can use this information to adjust the selection criteria, which may include adjusting the `CORPUS_MIN_INFO_GAIN` and using the `CORPUS_MAX_INFO_GAIN` pass-throughs.

If fewer generic terms are desired, increase the value of `CORPUS_MIN_INFO_GAIN` in small increments (0.5 or 1.0). If more generic terms are desired, decrease this value. `CORPUS_DEBUG` can be used to select a particular value of `CORPUS_MIN_INFO_GAIN`.

### **CORPUS\_DEBUG**

Recommendation: Set this pass-through to `true` only when you are tuning the filtering parameters; otherwise, do not use it.

## **Record-level filtering best practices**

### **RECORD\_NTERMS and RECORD\_FRACT\_OF\_MEDIAN**

Recommendation: The use of these pass-throughs depends on the length of text in the text property that contains candidate terms. The two scenarios considered here are properties with either short text or long text.

For short text (such the `P_Description` property in the wine data set shipped with the sample reference implementation), the recommendation is to not use these pass-throughs, which will keep all the terms.

For long text (such as news sites or sites with long articles), use the range version of `RECORD_NTERMS` to set however many terms per record you want, say a range of 16-24 or, if more is wanted, a range of 24-30. (Keep in mind that the lower limit should be greater than 10 and the upper limit should not be much larger than the lower limit.) Set `RECORD_FRACT_OF_MEDIAN` to 1.1 for relatively small documents, 1.2 for larger documents, and 1.5 for very large documents.

## **Format of the source data**

Terms can be extracted from either structured and unstructured source data.

In general, data and content acquisition will typically be used to retrieve the records used by the term extractor. In particular, crawling (through the Oracle Commerce Web Crawler or the CAS Server) is a viable source of content for term extraction.

The input must be as clean as possible and as similar to natural language as the original data allows. The following list provides some recommendations about how to pre-process unstructured text that is fed to the term extractor:

- Remove anything from the property that is sent to the term extractor that is not the main contents of the document. For example, when dealing with news articles, it is a good idea to remove bylines, copyright disclaimers, and so forth. When dealing with Web pages, the task is noticeably harder, because the navigational elements, page headers and footers, guestbook links, ads, etc., all have to be removed. In such extreme cases, one suggestion is to retain long sequences of sentences with correct sentence-terminating punctuation that does not have many major HTML tags (H1, P, HR, DIV, SPAN, and so on.) embedded: meaningful text is likely to satisfy this requirement, and items such as menus, ads, and page elements are likely to fail it.
- Remove anything that is not natural language text. This includes HTML tags, other markup, and long sequences of non-alpha characters (e.g., long sequences of dashes used as delimiters). Links to images, URLs (that might be used in plain text, outside of HTML tags), and anything that is not in grammatically correct language should be stripped. The same caveat applies to sequence of characters that are too long to be meaningful terms. The term extractor will report and skip those overlong noun phrases. However, it is useful to detect these upstream of the term extractor because their presence might indicate sections of the documents that should be removed.
- Punctuation should be correctly spaced, especially when stripping HTML or adding sentence-terminating punctuation. A sentence terminator is correctly interpreted only if it is followed by a space. For example:

```
Look at this.<IMG a.gif>Or this.
```

should not be converted to:

```
Look at this.Or this.
```

but instead to:

```
Look at this. Or this.
```

- Convert non-sentence text into sentences. If, for example, a useful section of the document is written as a list or a table (that is, separated with <LI> or <TD> tags), it is recommended to terminate such entries with periods (or semicolons, depending on context), if they are not so terminated to begin with.
- Merge text fields, if needed. For example, if the document title is a separate property, it is useful to append it to the main text property (terminating it with punctuation, if possible).
- Use the correct case for capitalized and non-capitalized text. If capitalization is not available, it makes a best guess and this guess is better when dealing with lower-case text than with all upper-case text. If the document text is in all upper case, it is advisable to convert it to all lower-case (or, possibly, all lower case with initial capitalization for the first word in every sentence).
- Use correct spelling in Web pages, especially blogs. For documents written in informal language, it is recommended that simple pattern replacement be done on most frequent terms (for example, "u" should be replaced with "you").

Whenever possible, the text should be cleansed in the source data. However, you can add record manipulators to the pipeline to perform pre-processing cleanup.

You can also use the CORPUS\_REGEX\_KEEP and CORPUS\_REGEX\_SKIP pass-throughs in the CAS manipulator to control which extracted terms are kept or discarded. For details on how to construct regular expressions with these pass-throughs, see the *Using regular expressions* topic of this chapter.



## Chapter 3

# Configuration Guidelines for Clustering

This section describes some guidelines for configuring cluster discovery.

## Configuration for clusters

Cluster Discovery configuration is passed as a part of the MDEX Engine configuration from a CAS-based application.

You must create the refinement configuration at

`<EAC_APP>\config\mdex\<EAC_APP>.refinement_config.xml` and make an entry for the configured dimension.

In the following example of the refinement configuration file, `P_Terms` is the dimension that is mapped to the extracted terms property - `OUTPUT_PROP_NAME`. The `CLUSTERS` tag contains the clustering parameters that the MDEX Engine reads during indexing.

`<EndecaApp>.refinement_config.xml`

```
<?xml version="1.0" encoding="UTF-8" standalone="no" ?>
<!DOCTYPE REFINEMENTS_CONFIG SYSTEM "refinement_config.dtd">
<REFINEMENTS_CONFIG>

.....

<REFINEMENTS DVAL_COLLAPSE_THRESHOLD="" NAME="P_Terms" SORT_TYPE="ALPHA">
  <STATS NUM_RECORDS="FALSE"/>
  <DYNAMIC_RANKING COUNT="10" MORE="FALSE" TOP_REFINEMENTS_SORT="DEFAULT"
TYPE="FREQUENCY"/>
</REFINEMENTS>

  <CLUSTERS COHERENCE="8" MAX_CLUSTERS="10" MAX_CLUSTER_OVERLAP="5" MAX_CLUSTER_SIZE="5" MAX_REFINEMENT_PRECISION="0.250000" NAME="P_Terms" REC_SAMPLE_SIZE="500"/>
</CLUSTERS>
</REFINEMENTS_CONFIG>
```

The following topics provide descriptions of the parameters and some guidelines for tuning them.

## Clustering parameter descriptions

This topic describes clustering parameters.

The list below describes the clustering parameters, including their range of values with defaults and recommended values. The next topic in this section explains the rationale for the recommended values and provides an order in which these parameters should be adjusted.

### Sample size

Description: Clustering is done by examining term distribution across a sample of the result set. This parameter governs how many records are sampled from the navigation state. Clustering processing time and memory consumption are both roughly linear with this number; thus, lowering the value results in smaller memory consumption and faster turnaround. However, statistical errors are likely to occur when the sample size is small. Setting this value higher overcomes statistical errors for data sets where fewer terms are tagged onto each record.

Range: Integer, 50-2000 (default: 500)

Recommended value: 500

### Maximum clusters

Description: This parameter limits the number of clusters that are generated by the MDEX Engine.

Range: Integer, 2-10 (default: 10)

Recommended value: 6

### Coherence

Description: This parameter governs the decision of whether a set of terms is coherent enough to form a cluster (each cluster should have only closely related records). Low values are permissive (for example, not demanding much coherence) and results in fewer larger clusters. High values are strict and result in smaller clusters. An average value is recommended.

Range: Integer, 0-10 (default: 5)

Recommended value: 5

### Maximum precision

Description: Terms that are extracted from sampled records are filtered by their precision **p** where **p** = number of sampled records that this term is tagged onto divided by the number of all sampled records. Terms that have too high a value of **p** are likely to be the search term or be synonymous with it or be too general to make for a good clustering term. If you use the recommended tuning values of the term extractor, each term is tagged to only roughly 1/3 of the records that contain this term in the text, which means that the search term, if present, has **p** of roughly 0.33 - more or less stringent tuning of the term extractor changes this value. There usually is a gap in the values of **p** between the search term and the more useful terms, which start at approximately **p** = 0.25 and less.

Range: Float, 0.0 - 1.0 (default: 1.0)

Recommended value: 0.25

### Maximum cluster size

Description: This parameter sets the maximum number of terms in a cluster. Each cluster has at least 2 terms. Because of the match-partial cluster selection mechanism, the more terms there are in the cluster, the potentially



higher its coverage is. On the other hand, the clusters that are too large take up too much space to display and take too long for users to read.

Range: Integer, 2 - 10 (default: 10)

Recommended value: 8

### Maximum cluster overlap

Description: If two clusters overlap (the record sets that each cluster maps to overlap), then the smaller one, as measured by the estimated size of the record set it maps to, can be removed, depending on how big this overlap is. This parameter dictates the overlap above which the smaller cluster is removed.

Clusters which overlap by more than this value are removed. Thus, the default setting of 10 means that clusters that overlap by more than 10 out of 10 records are removed. Since this is impossible, this means that setting of 10 disables cluster overlap filtering, which is most extreme level of coarseness for this filter. Tuning this parameter down makes the cluster overlap more and more fine-grained. Thus, a value of 9 removes only the clusters that greatly overlap; setting it to the recommended value of 5 removes only clusters overlapping half-way or so (remember that the overlap is merely estimated). Setting this parameter to lower values (less than 5) makes overlap filtering quite sensitive and removes clusters which overlap even by a small amount. Note that clusters that do not overlap are never filtered.

Range: Integer, 0-10 (default: 10)

Recommended value: 5

## Tuning strategy for clusters

This topic provides guidelines for tuning the clustering parameters.

The guidelines for the clusters tuning strategy include initial values that are used for the first trial clustering run and recommended values. The tuning process involves changing the parameters from their initial values to their recommended values, with certain variation dependent on the properties of the particular data set and the application needs.

In general, the tuning strategy involves starting with the parameters at a permissive setting and then gradually decreasing the value. You tune the parameters by observing their impact simultaneously on the results for several different queries: no query or node 0, broad queries, narrow queries, single-term query, and multi-term query. In other words, you should avoid tuning the parameters based on a specific query.

The following procedure is intended as a tool for gradual tuning. It allows you to observe the effect of changing the parameters on several different queries at once. Use the suggested order. It maps to the order in which these parameters impact the clustering algorithm, from upstream to downstream.

### 1: Number of records sampled from the navigation state

Recommended value: 500

Initial value: 500

Strategy: Start with the parameter set to 500, and increase it if you see that the terms at the bottom of your related terms list (terms 100-120 or so) are seen in fewer than 3 records.

Note that the recommended value of 500 is for data sets with 20 or more terms tagged onto each record. Use a higher value for data sets with fewer terms per record. If an average record has only 2 to 3 terms per record, set this value to 2000. A good rule of thumb for this value is: when the 120 most frequent terms are sampled

for clustering, the 120th term should be present in at least 3 records. If it is present in fewer, this setting should be increased.

## **2: Maximum refinement precision**

Recommended value: 0.25

Initial value: 1.0

Strategy: Start with this value set to 1.0 (no precision filtering). Try several different queries and pick a level of top useful precision that separates useful terms from the frequent but uninformative ones. Note that, typically, only the values between 0.05 and 0.5 are useful.

## **3: Maximum number of terms per cluster**

Recommended value: 6 - 8

Initial value: 10

Strategy: Start with a value of 10 to see all the terms that are getting into clusters. Reduce the value until the clusters are small enough to fit into whatever real estate your UI provides. Using a value of 2 is not recommended. Note that the cluster coverage and recall are reduced when the number of terms is reduced.

## **4: Cluster Coherence**

Recommended value: 5

Initial value: 5

Strategy: Start with the default value of 5. If you see undesirable cluster splintering - several clusters that seem to map to the same semantic areas - then this value should be decreased. On the other hand, if the cluster set is missing some semantic areas, this value should be increased. Note that it is acceptable to have several overlapping clusters remaining after tuning this value, because they are removed in the next step.

## **5: Maximum cluster overlap**

Recommended value: 5

Initial value: 10

Strategy: Start with a value of 10, then decrease this parameter until the desired number of overlapping clusters remains. In some cases, for example, depending on customer needs, some cluster overlap can be retained, particularly if the smaller cluster is an especially coherent one).

## **6: Maximum number of clusters**

Recommended value: 6

Initial value: 10

Strategy: Start with a value of 10 to see all the available clusters after all the other settings had been applied. Reduce this number if you still see more clusters than permitted by the available UI space.

## Building the Front End of the Term Discovery Application

This section provides information on how to change your front-end application to display terms and clusters.

### Files to be changed

You can modify existing JSP files and add a new file.

The examples are based on the following JSP files that are shipped with the `endeca_jspref` reference implementation:

- `constants.jsp`
- `controller.jsp`
- `nav_controls.jsp`

In addition, a new `nav_clusters.jsp` file is provided in Appendix A as an example of rendering clusters.

### Adding global constants

The `constants.jsp` file sets global variables that can be used in any other page in the application.

It is a good practice to add global variables for handling the displaying of the terms dimension and the clusters. Doing so will allow you to easily change display characteristics in one central file.

The following Java code is an example of a section that can be added to the `constants.jsp` file:

```
// Dimension name of the Term Discovery dimension.
// Make it null if you do NOT want term processing.
private static final String relTermsDimName = "Terms";

//Display name of the Term Discovery dimension.
private static final String relTermsDisplayString = "Term Discovery";

//The rootId of the Term Discovery dimension.
//Make it -1 if you do NOT want terms processing.
private static final long relTermsRootId = 21;

// Handling for the all-terms property.
// Ignored if String is null.
```

```
private static final String P_AllTerms = "P_AllTerms";

// if true, a More... link shows after the term's short list.
// Should be false in production.
private static final boolean showTermsMore = false;

// If true, TD dimension and property show in record display
// Should be false in production.
private static final boolean showTermsInRecord = false;
```

The code sets up the following variables. These constants will be used in most of the JSP pages listed above.

Global Variable	Purpose
relTermsDimName	Sets the name of the Guided Search dimension for Term Discovery.
relTermsDisplayString	Sets the name that will be displayed in application pages for the Term Discovery dimension.
relTermsRootId	Sets the root ID of the Term Discovery dimension.
P_AllTerms	Sets the name of the Guided Search property that contains all the terms.
showRelTermsMore	Sets whether a More... link is shown for terms.

## Setting refinements in the controller file

The `controller.jsp` module is the entry point into the Guided Search application.

The controller file receives the browser request from the application server, formulates the query, and sends the query to the MDEX Engine.

The following code should be added to the `controller.jsp` file:

```
if (relTermsRootId >= 0 && usq.containsNavQuery()
    && request.getParameter("Ne") == null) {
    DimValIdList dvl = new DimValIdList();
    dvl.addDimValueId(relTermsRootId);
    usq.setNavExposedRefinements(dvl);
}
```

The code can be added after the `ENEQuery` object named `usq` is created.

The code first tests that three conditions are true:

- The `relTermsRootId` global variable is set to a value greater than 0.
- The `ENEQuery.containsNavQuery()` method determines that the current query is a navigation query.
- The request does not include the `Ne` URL parameter, which determines which dimension navigation refinements are exposed.

If all three conditions are true, then the code performs three actions:

1. Creates an empty `DimValIdList` object.
2. Uses the `DimValIdList.addDimValueId()` method to add the ID of the Term Discovery dimension to the list.
3. Uses the `ENEQuery.setNavExposedRefinements()` method to set the refinements that are exposed/returned with the navigation query.

In step 3, the dimension ID of the Term Discovery dimension is given as the argument, which means that this dimension will be the parent of the refinements that should be returned with the query. The returned refinements will be the terms.

## Displaying refinements

The terms in a Term Discovery dimension are returned as dimension value refinements.

These refinements are displayed according to the coding in the `nav_controls.jsp` file. Because a number of changes are required to handle the Term Discovery refinements, the modified file is included in Appendix A.

## Displaying clusters

The information for each cluster is returned from the MDEX Engine in a `Supplement` object that accompanies the result of a navigation query.

The `Supplement` objects are encapsulated in a `SupplementList` object.

## Cluster properties

Each cluster (`Supplement` object) contains a `PropertyMap` object.

The `PropertyMap` object in turn contains the following cluster-related properties (as key/value pairs):

Key Name	Value
Dgraph.SeeAlsoCluster	The name of the Term Discovery dimension from which this cluster was generated.
ClusterRank	The zero-based rank of this cluster (0 for the first cluster, 1 for the second cluster, and so on).
NTerms	A number indicating the number of terms in this cluster.
Term_i (where i is 0, 1 ... NTerms-1)	The term at rank i within this cluster. Note that there will be as many Term_i entries as there are terms in the cluster.

## JSP code for displaying clusters

You can add a new `nav_clusters.jsp` file for rendering clusters.

Appendix A of this guide includes a new `nav_clusters.jsp` file that can be used as a template for rendering cluster contents. This file should be included at the end of the `nav_controls.jsp` file. Note that it is recommended that you display clusters only if there are two or more.

The highlights of this file are as follows:

1. Get all the `Supplement` objects from the `Navigation` object, which will be encapsulated in a `SupplementList` object (note that the list can also include non-cluster `Supplement` objects):

```
SupplementList navsups = nav.getSupplements();
```

2. Within a For loop (labeled supLoop), get a Supplement object and then get that object's properties:

```
Supplement sup = (Supplement)navsups.get(I);
PropertyMap propsMap = sup.getProperties();
```

3. Get the value of the Dgraph.SeeAlsoCluster property:

```
String clustersPropName = (String)propsMap.get("DGraph.SeeAlsoCluster");
```

4. Test the value returned from the Dgraph.SeeAlsoCluster property. If it is null, then this Supplement object is not a cluster and we should loop back to step 2; if the value is non-null, then this is a cluster and we continue to step 5:

```
if (clustersPropName != null)
```

5. If this is the first discovered cluster, then the clustersOn variable will be set to false. Therefore, first display the Cluster Discovery header and then set the clustersOn variable to true (so that the header will not be displayed again):

```
if (!clustersOn) {
    // display title
    ...
    clustersOn = true;
}
```

6. Get the ranking (ClusterRank property) of the cluster and the number of terms (NTerms property) it contains. The returned string values are then transformed to integers:

```
String rankString = (String)propsMap.get("ClusterRank");
String nTermsString = (String)propsMap.get("NTerms");
int rank;
int nTerms;
try {
    rank = Integer.parseInt(rankString);
    nTerms = Integer.parseInt(nTermsString);
}
```

7. Within a For loop, retrieve the terms from the NTerms property and format them by separating them with commas and spaces. The list of terms will then be:

```
StringBuffer termsSB = new StringBuffer();
StringBuffer termsSBSpace = new StringBuffer();
for (int iTerm = 0; iTerm < nTerms; ++iTerm) {
    String term = (String)propsMap.get("Term_"+iTerm);
    ...
    if (termsSB.length() != 0) {
        termsSB.append(", ");
        termsSBSpace.append(" ");
    }
    termsSB.append(term);
    termsSBSpace.append(' ').append(term).append(' ');
}
```

8. For a given cluster, begin to create a UrlENEQuery request (using the current request), in case the user wants to click on that cluster. Also get the current navigation searches (as an ERecSearchList) to determine if the cluster selection is already active in the searches.

```
UrlENEQuery newq = new UrlENEQuery(request.getQueryString(),"UTF-8");
ERecSearchList searches = newq.getNavERecSearches();
```

9. Create a new search (an `ERecSearch` object), using the `clusterPartial` search interface as the search key, the list of related terms (in the `clusterSpace` variable) as the terms for the search, and mode `matchpartial` as the search option (which specifies `MatchPartial` as the search mode).

```
ERecSearch newSearch = new ERecSearch("clustersPartial",
    clusterSpace, "mode matchpartial");
```

10. Test whether the current navigation searches are null or do not contain the new search (from step 9). If the test is true, then the new search can be added to the `UrlENEQuery` request; if it is false, do not add the new search because the cluster selection is already active.

```
if (searches == null || !searches.contains(newSearch)) {
    ...
    searches.add(newSearch);
    newq.setNavERecSearches(searches);
    ...
}
```

11. Loop back to step 2 to get another `Supplement` object. The loop is done when all the objects in the `SupplementList` have been retrieved.

12. Display the clusters, which are stored as a list of strings in the `clusterStrings` variable. The `clusterUrls` variable is a list of the cluster URLs:

```
for (int I = 0; I < clusterStrings.size(); ++I) {
    %><tr><td></td>
    <td width="100%"><font face="arial" size="1" color="gray">
    > <a href="<%= clusterUrls.get(I) %>">
    <font face="arial" size="1" color="blue">
    <%= clusterStrings.get(I) %></font></a>
    </td></tr><%
}
```

The following is an abbreviated example of the JSP reference implementation showing the clusters rendered by the `nav_clusters.jsp` file. Clicking on a cluster link will execute the partial match query built by steps 8-10.



## Clustering overlap properties

Clustering overlap information is also returned by the MDEX Engine.

The `PropertyMap` object (in the cluster `Supplement`) also includes the following set of properties that provide clustering overlap information.

Key Name	Value
<code>Dgraph.SeeAlsoClusterOverlaps</code>	The name of the Term Discovery dimension from which this cluster was generated.
<code>NClusters</code>	A number indicating the number of clusters that were returned by the MDEX Engine.
<code>Cluster_i</code> (where <code>i</code> is 0, 1, ... <code>NClusters-1</code> )	The cluster overlap numbers for a given cluster. Note that the cluster number (the <code>i</code> value) corresponds to the <code>ClusterRank</code> value in the <code>DGraph.SeeAlsoCluster</code> object.

These properties provide a square matrix that has the cluster overlap numbers. In the matrix, `number (i, j)` is the estimated number of records (from the records sampled from the navigation states) that are covered by both cluster `i` and cluster `j`.

Note that from the definition it follows that diagonal numbers `(i, i)` have the estimated number of records covered by each particular cluster. These diagonal numbers tend to decrease, because of the way that the Cluster Discovery software sorts clusters (by decreasing estimated coverage).

This information can be used in application-specific ways, for example, by an application page that presents a graphical depiction of the clusters.

## Displaying records and dimension refinements

Records and refinements from a Term Discovery dimension are displayed like other dimensions.

There is no difference in displaying refinement dimension values from a Term Discovery dimension than from regular dimensions. Information about displaying refinements is found in the *MDEX Engine Development Guide*, in the chapter titled “Working with Dimensions.”

Likewise, the process of displaying Guided Search records generated from Term Discovery refinements is the same as with any Guided Search record. For details, see the chapter titled “Working with Endeca Records” in the *MDEX Engine Development Guide*.



## Chapter 5

# Term Discovery Advanced Topics

This section discusses advanced topics for Term Discovery applications.

## Term filtering with pre-tagged records

The Term Discovery software can apply filtering to documents that have already been tagged.

The use case described in this scenario involves a corpus in which the documents are tagged with pre-existing terms that were generated and maintained by an external process -- not by the Guided Search term extraction software. The goal is to apply corpus-level and record-level filtering to these terms, just as though they had been identified as candidates by the term extractor itself.

There are three variations of this use case:

- Filtering only of pre-tagged terms: no new terms are extracted from the records, but you want filtering to be performed on the pre-tagged terms.
- Uniform filtering of both sets of terms: the term extractor extracts new terms and combines them with the pre-tagged terms. The same filters are applied to both sets of terms equally.
- Filtering only one of the sets of terms: new terms are extracted from the records by the term extractor. Only the newly-extracted terms are filtered, but the pre-tagged terms are not. After filtering, both sets of terms are combined into one output property and tagged on the records.



**Note:** In all cases, the pre-tagged source property must contain only one term as its value. Any given record can have multiple instances of this property.

## Filtering only pre-existing terms

This use case assumes that the source records have already been tagged with terms.

In this use case, you want to perform corpus- and/or record-level filtering on these pre-tagged terms. However, you do not want the term extractor to extract any more terms from the records.

The CAS manipulator that will perform filtering on the pre-tagged terms should have the following configuration values for the pass-throughs:

PASS_THROUGH Element	Configuration Value
RECORD_SPEC_PROP_NAME	Set to the name of the record specifier property.

<b>PASS_THROUGH Element</b>	<b>Configuration Value</b>
TEXT_PROP_NAME	Set to the name of a non-existent property, so that no new terms are extracted.
INPUT_TERM_PROP_NAME	Set it to the name of the source property containing the pre-tagged terms.
OUTPUT_PROP_NAME	Set it to the name of the property that is the destination for the pre-existing terms on the Guided Search record.
corpus-level pass-throughs	As required by the application.
record-level pass-throughs	As required by the application.

## Filtering both sets of terms uniformly

This use case assumes that you want to perform term extraction on a data set that already has pre-tagged terms.

In this use case, you want to combine both sets of terms and have the same filtering applied to them.

The CAS manipulator that will perform uniform filtering on both sets of terms should have the following configuration values for the pass-throughs:

<b>PASS_THROUGH Element</b>	<b>Configuration Value</b>
RECORD_SPEC_PROP_NAME	Set to the name of the record specifier property.
TEXT_PROP_NAME	Set to the name of the source text property from which new terms will be extracted.
INPUT_TERM_PROP_NAME	Set to the name of the source property containing the pre-tagged terms.
OUTPUT_PROP_NAME	Set to the name of the property that is the destination for both newly-tagged terms and pre-existing terms.
corpus-level pass-throughs	As required by the application.
record-level pass-throughs	As required by the application.

When the term extractor runs, both newly-extracted terms and pre-tagged terms are output to the same OUTPUT\_PROP\_NAME property. As a result, the same corpus- and record-filtering is applied to all terms.

## Filtering only the new terms

This use case assumes that new terms are extracted from the records by the term extractor, but are not immediately combined with the pre-tagged terms.

In this use case, the newly-extracted terms are not combined at first with the pre-tagged terms. Instead, only the newly-extracted terms are filtered, and the pre-tagged terms are not. After filtering, both sets of terms are combined into one output property and tagged on the Guided Search records.

The CAS manipulator that will perform this type of filtering should have the following configuration values for the pass-throughs:

<b>PASS_THROUGH Element</b>	<b>Configuration Value</b>
RECORD_SPEC_PROP_NAME	Set to the name of the record specifier property.

PASS_THROUGH Element	Configuration Value
TEXT_PROP_NAME	Set to the name of the source text property from which new terms will be extracted.
OUTPUT_PROP_NAME	Set to the name of the source property containing the pre-tagged terms. That is, the destination for the new terms will be the same property as the pre-existing terms.
corpus-level pass-throughs	As required by the application.
record-level pass-throughs	As required by the application.

Note that the INPUT\_TERM\_PROP\_NAME pass-through is not used. As a result, the pre-existing terms are not filtered, but are used as is.

When the term extractor runs, only the newly-extracted terms are filtered with the corpus- and record-level pass-throughs. The filtered terms are then output to the OUTPUT\_PROP\_NAME property, which is the name of the property with the pre-existing terms. As a result, the same corpus- and record-filtering is applied to all terms. Note that if duplicate values are created for the OUTPUT\_PROP\_NAME property, they are removed by the property mapper.

## Tuning aids for the filtering parameters

This section discusses two tuning aids that will help you when you are tuning the parameters for the corpus- and record-level pass-throughs.

The two tuning aids are:

- STATEFUL update mode
- Corpus-verbose logs

### Using STATEFUL mode for tuning

You can use STATEFUL mode to tune the values you set for the corpus-level and record-level pass-throughs in the CAS manipulator.

Before you begin, make sure the baseline pipeline has the UPDATE\_MODE pass-through set to STATEFUL mode.

The general procedure for using STATEFUL mode for tuning is:

1. Generate a baseline update, index the records, and start the MDEX Engine.
2. Run searches against the P\_AllTerms property and check the quality of the clusters.
3. Adjust the corpus-level and/or record-level parameters.
4. Add a MAX\_INPUT\_RECORDS pass-through set to 0 (zero) to the CAS manipulator.
5. Generate another baseline update. The update will be much faster because no terms will be extracted. However, a full corpus- and record-level filtering operation will be performed.
6. Repeat the above steps (except step 4) until you are satisfied with the results.

When you finish, be sure to remove the MAX\_INPUT\_RECORDS pass-through so all your source records will be processed.

## Using corpus-filtering logging statistics

The CORPUS\_DEBUG pass-through is helpful for generating debugging information.

The CORPUS\_DEBUG pass-through enables the term extractor to log detailed information about the scores that it assigns to terms. Temporarily setting this pass-through helps you tune corpus-level filtering.

The log entries contain four fields of information:

Log Entry	Information
term	The term (noun phrase) that was extracted and filtered.
count	The number of documents in which this term occurs at least once. You can use the RECORD_NTERMS pass-through to set a limit on the number of documents in which a term can occur.
coverage	The coverage score, which is a percentage of all the corpus documents in which this term was found. You can use the CORPUS_MIN_COVERAGE and CORPUS_MAX_COVERAGE pass-throughs to adjust the percentage.
info_gain	the info_gain score, which is a measure of the global informativeness of the term. The CORPUS_MIN_INFO_GAIN and CORPUS_MAX_INFO_GAIN pass-throughs will affect this score.

An example of a term log entry is:

```
term: airport count: 60 coverage: 0.04 info_gain: 1.614589
```

In this example, the term `airport` was found in 60 documents, which is 0.04 (4%) of the 1500-document corpus, and a global informativeness score of 1.614589 was given to the term.

## Term Discovery Sample Files

## Modified nav\_controls.jsp file

The following `nav_controls.jsp` file has been modified to display refinements from the Terms Discovery dimension. Added or modified code is highlighted in bold face.

```

-----<%
DESCRIPTION:
This module displays basic, standard navigation controls.  It
is mainly used for debugging purposes and as a starting point for
no-frills navigation solutions. It only displays refinement dimensions,
so this module should be used in conjunction with nav_breadcrumbs.

Copyright (C) 2008 by Endeca Technologies - COMPANY CONFIDENTIAL
----->%

<table border="0" cellspacing="0" cellpadding="0" width="100%">
  <tr><td colspan="2" bgcolor="orange"><font face="arial" size=2 color="white">
    &nbsp; nav_controls:<font></td></tr>
  <tr><td colspan="2"></td></tr>
  <%
    // Get refinement dimension groups
    DimGroupList refDimensionGroups = nav.getRefinementDimGroups();
    // Get descriptor dimensions
    DimensionList descDimensionsNC = nav.getDescriptorDimensions();
    // Output message if no refinement options left
    if (refDimensionGroups.size() == 0) {
      %>
      <tr><td colspan="2"><font face="arial" size=3 color="orange">
        <i>No Additional Query<br>Parameters Available</i></font></td></tr>
      <%
    }
    // Output message if no refinement options have been made
    else if (descDimensionsNC.size() == 0) {
      %>
      <tr><td colspan="2"><font face="arial" size=3 color="orange">
        <i>Query Parameters:</i></font></td></tr>
      <tr><td colspan="2"></td></tr>
      <%

```

```

}
// Header if additional query parameters available
else {
    %>
    <tr><td colspan="2"><font face="arial" size=3 color="orange">
    <i>Additional Query Parameters:</i></font></td></tr>
    <tr><td colspan="2"></td></tr>
    <%
}
// Loop over dimension groups
for (int i=0; i<refDimensionGroups.size(); i++) {
    // Get dimension group object
    DimGroup dg = (DimGroup)refDimensionGroups.get(i);
    // If group is explicit (not default group), display group
    if (dg.isExplicit()) {
        %>
        <tr><td colspan="2"></td></tr>
        <tr><td colspan="2"><font face="arial" size="2" color="#999999">
        <%= dg.getName() %></font></td></tr>
        <%
    }
    // Loop over dimensions in group
    for (int j=0; j<dg.size(); j++) {
        // Get dimension object
        Dimension dim = (Dimension)dg.get(j);
        // Get root for dimension
        DimVal root = dim.getRoot();
        // Get id of root
        long rootId = root.getId();
        // special handling for Term Discovery dimension
        final boolean isRelTerms = dim.getName().equals(relTermsDimName);
        // Get refinement list for dimension
        DimValList refs = dim.getRefinements();
        // Create request to expose dimension values
        UrlGen urlg = new UrlGen(request.getQueryString(), "UTF-8");
        urlg.removeParam("D");
        urlg.removeParam("Dx");
        urlg.removeParam("sid");
        urlg.removeParam("in_dym");
        urlg.removeParam("in_dim_search");
        urlg.addParam("sid", (String)request.getAttribute("sid"));
        // Expand dimension
        if (!isRelTerms && refs.size() == 0) {
            urlg.addParam("Ne", Long.toString(rootId)+
            (relTermsRootId>= 0? " "+Long.toString(relTermsRootId):""));
        }
        // Close dimension
        else {
            urlg.removeParam("Ne");
        }
        String url = CONTROLLER+"?" +urlg;
        // Display dimension (open row here, close later)
        if (!isRelTerms) {
            %>
            <tr><td colspan="2"><a href="<%= url %>">
            <font face="arial" size="2" color="#444444">
            <%= dim.getName() %></font></a>
            <%
        }
        else {
            %>

```

```

<tr><td colspan="2" bgcolor="orange">
<font face="arial" size=2 color="white">
&nbsp;  <%=relTermsDisplayString %><font></td></tr>
<%
}

// Get intermediate list for dimension
DimValList ints = dim.getIntermediates();
// Loop over intermediate list
for (int k=0; k < ints.size(); k++) {
    // Get intermediate dimension value
    DimVal intermediate = ints.getDimValue(k);
    // Display intermediate
    %><font face="arial" size="2" color="#444444"> >
    <%= intermediate.getName() %></font><%
}
// Close nav row
%></td></tr><%
String refinementsColor = "blue";
Set activeDiscTerms = new HashSet();
if (isRelTerms) {
    String ntk = (String)request.getParameter("Ntk");
    String ntx = (String)request.getParameter("Ntx");
    if (ntk != null && ntk.equals(P_AllTerms) &&
        "mode matchall".equals(ntx))
    {
        String discTerm = (String)request.getParameter("Ntt");
        if (discTerm != null) {
            if (discTerm.length() >= 3 &&
                discTerm.charAt(0) == '"' &&
                discTerm.charAt(discTerm.length()-1) == '"')
            {
                // remove quotes
                discTerm = discTerm.substring(1, discTerm.length()-1);
            }
            activeDiscTerms.add(discTerm);
        }
    }
} // if (isRelTerms)

%><%
// Loop over refinement list
for (int k=0; k < refs.size(); k++) {
    // Get refinement dimension value
    DimVal ref = refs.getDimValue(k);
    // Get properties for refinement value
    PropertyMap pmap = ref.getProperties();
    // Get dynamic stats
    String dstats = "";
    if (pmap.get("DGraph.Bins") != null) {
        dstats = " (" + pmap.get("DGraph.Bins") + ")";
    }
}
%><%
// Create request to select refinement value
urlg = new UrlGen(request.getQueryString(), "UTF-8");
boolean displayRefinement=true;
if (isRelTerms &&
    (!ref.getName().equals("More...") || !showRelTermsMore))
{
    if (ref.getName().equals("More...") ||
        ENEQueryToolkit.isImplicitRefinement(dim, ref) ||

```

```

        activeDiscTerms.contains(ref.getName()))
    {
        displayRefinement = false;
    }
    else {
        urlg.addParam("Ntk","clustersPartial");
        urlg.addParam("Ntt","\\"+ref.getName()+"\\");
        urlg.addParam("Ntx","mode matchall");
        urlg.removeParam("No");
        urlg.removeParam("Nao");
        urlg.removeParam("Nty");
        urlg.removeParam("D");
        urlg.removeParam("Dx");
        urlg.removeParam("sid");
        urlg.removeParam("in_dym");
        urlg.removeParam("in_dim_search");
        urlg.addParam("sid",(String)request.getAttribute("sid"));
        url = CONTROLLER+"?" +urlg;
    }
%><%
} else {
    // If refinement is navigable, change the Navigation parameter
    if (ref.isNavigable()) {
        urlg.addParam("N",
            (ENEQueryToolkit.selectRefinement(nav,ref)).toString());
        urlg.addParam("Ne",Long.toString(rootId)+
            (relTermsRootId>= 0? " "+Long.toString(relTermsRootId):""));
    }
    // If refinement is non-navigable, change only the
    // exposed dimension parameter
    // (Leave the Navigation parameter as is)
    else {
        urlg.addParam("Ne",Long.toString(ref.getId()+
            (relTermsRootId>= 0? " "+Long.toString(relTermsRootId):"")));
    }
    urlg.removeParam("No");
    urlg.removeParam("Nao");
    urlg.removeParam("Nty");
    urlg.removeParam("D");
    urlg.removeParam("Dx");
    urlg.removeParam("sid");
    urlg.removeParam("in_dym");
    urlg.removeParam("in_dim_search");
    urlg.addParam("sid",(String)request.getAttribute("sid"));
    url = CONTROLLER+"?" +urlg;
}
// Display refinement
if (displayRefinement) {
    %>
    <tr><td></td>
    <td width="100%"><a href="<%= url %>">
    <font face="arial" size="1" color="<%=refinementsColor%>">
    <%= ref.getName() %></font></a>
    <font face="arial" size="1" color="gray"><%= dstats %>
    </font></td></tr>
    <%
    }
    } // end of: Loop over refinement list
} // end of: Loop over dimensions in group
// If group is explicit (not default group), display spacer
if (dg.isExplicit()) {

```



```

        %>
        <tr><td colspan="2"></td></tr>
        <%
    }

    } // end of: Loop over dimension groups
    %>
    <tr><td colspan="2"></td></tr>
</table>

<%-- Display Clusters Controls --%>
<%@ include file="nav_clusters.jsp" %>

<%-- Display Range Filter Controls --%>
<%@ include file="nav_range_controls.jsp" %>

```

## New nav\_clusters.jsp file

This nav\_clusters.jsp sample file is used to render clusters that are generated by the Cluster Discovery feature.

This file should be included in the nav\_controls.jsp file.

```

<%-----
DESCRIPTION:
This module demonstrates the Cluster Discovery feature.
It displays clusters received as Supplemental Objects, makes them
selectable, and, upon selection, generates a search based on
the selected clusters.

This module is included in the nav_controls.jsp module.

Copyright (C) 2008 by Endeca Technologies - COMPANY CONFIDENTIAL
-----%>
<%
// Get supplemental list
SupplementList navsups = nav.getSupplements();
boolean clustersOn = false;

// lazily allocated:
List<String> clusterUrls = null;
List<String> clusterStrings = null;

// Loop over cluster supplemental objects
supLoop:
for (int i = 0; i < navsups.size(); ++i) {
    // Get individual see also object
    Supplement sup = (Supplement)navsups.get(i);
    // Get property map
    PropertyMap propsMap = sup.getProperties();

    String clustersPropName = (String)propsMap.get("DGraph.SeeAlsoCluster");
    if (clustersPropName != null) {
        if (!clustersOn) {
            // display title
            %>
            <table border="0" cellspacing="0" cellpadding="0" width="100%">
            <tr><td colspan="5" bgcolor="orange">

```

```

        <font face="arial" size=2 color="white">
Cluster Discovery</font></td></tr>
<tr><td colspan="5"></td></tr>
<%
clustersOn = true;
} // end of if !clustersOn
String rankString = (String)propsMap.get("ClusterRank");
String nTermsString = (String)propsMap.get("NTerms");
int rank;
int nTerms;
try {
    rank = Integer.parseInt(rankString);
    nTerms = Integer.parseInt(nTermsString);
} catch (NumberFormatException e) {
    // add code here to log error
    continue supLoop;
} // end of catch

StringBuffer termsSB = new StringBuffer();
StringBuffer termsSBSpace = new StringBuffer();
for (int iTerm = 0; iTerm < nTerms; ++iTerm) {
    String term = (String)propsMap.get("Term_"+iTerm);
    if (term == null) {
        // add code to log error
        continue supLoop;
    }
    if (termsSB.length() != 0) {
        termsSB.append(", ");
        termsSBSpace.append(" ");
    } // end of if termsSB.length
    termsSB.append(term);
    termsSBSpace.append(' ').append(term).append(' ');
} // end of for terms
String clusterX = termsSB.toString();
String clusterSpace = termsSBSpace.toString();
// Create request to follow a cluster selection
// (unless cluster selection already active in the searches)
UrlENEQuery newq = new UrlENEQuery(request.getQueryString(), "UTF-8");
ERecSearchList searches = newq.getNavERecSearches();
ERecSearch newSearch = new ERecSearch("All", clusterSpace,
    "mode matchpartial");

if (searches == null || !searches.contains(newSearch)) {
    if (searches == null)
        searches = new ERecSearchList();
    if (clusterUrls == null) {
        clusterUrls = new ArrayList<String>();
        clusterStrings = new ArrayList<String>();
    } // end of if clusterUrls
    searches.add(newSearch);
    newq.setNavERecSearches(searches);
    UrlGen hostportq = new UrlGen("", "UTF-8");
    hostportq.addParam("eneHost",
        (String)request.getAttribute("eneHost"));
    hostportq.addParam("enePort",
        (String)request.getAttribute("enePort"));
    String url = CONTROLLER + "?" + hostportq.toString() +
        "&" + UrlENEQuery.toQueryString(newq, "UTF-8");
    clusterUrls.add(url);
    clusterStrings.add(clusterX);
} // end of if searches

```

```

    } // end of if clusterPropName != null
} // end of if clusterPropName != null

if (clusterStrings != null && clusterStrings.size() > 1) {
    // display clusters only if at least 2.
    for (int i = 0; i < clusterStrings.size(); ++i) {
        %>
        <tr><td></td>
        <td width="100%"><font face="arial" size="1" color="gray">
        > <a href="<%= clusterUrls.get(i) %>">
        <font face="arial" size="1" color="blue">
        <%= clusterStrings.get(i) %></font></a>
        </td></tr>
        <%
    } // end of for clusterStrings
} // end of if clusterStrings

if (clustersOn) {
    // close table
    %>
    <tr><td colspan="2"></td></tr>
    </table>
    <%
}%> // end of if clusterOn

```



# Index

## A

ALL\_TERMS\_OUTPUT\_PROP\_NAME pass-through  
definition 17

## B

baseline updates, mode for 20  
best practices for term filtering 27

## C

clusters  
configuration parameters 32  
configuring 31  
JSP code for rendering 37  
overlap properties 40  
overview 11  
properties in Supplement objects 37  
tuning strategy 33  
constants.jsp file  
adding global constants 35  
setting refinements 36  
CORPUS\_DEBUG pass-through  
used for tuning 44  
CORPUS\_MAX\_COVERAGE pass-through  
recommended setting 28  
CORPUS\_MAX\_INFO\_GAIN pass-through  
recommended setting 28  
CORPUS\_MAX\_RECS pass-through  
recommended setting 27  
CORPUS\_MIN\_COVERAGE pass-through  
definition 24  
recommended setting 27  
CORPUS\_MIN\_INFO\_GAIN pass-through  
definition 24  
recommended setting 28  
CORPUS\_MIN\_RECS pass-through  
definition 23  
recommended setting 27  
CORPUS\_REGEX\_KEEP pass-through  
definition 24  
recommended setting 28  
CORPUS\_REGEX\_SKIP pass-through  
definition 24  
recommended setting 28  
corpus-level filtering  
best practices 27  
configuration parameters 23  
for new terms only 42  
for pre-tagged and extracted terms 42  
for pre-tagged terms 41  
coverage score for terms 44

## D

destination property for tagged terms 17

## E

exclude list for term extraction 19

## F

filtering applied to pre-tagged records 41  
format of source data 29

## G

global constants for the front-end application 35  
global informativeness of terms, threshold for 24  
global language ID for documents, setting 22  
Guided Search Cluster Discovery, overview of 11  
Guided Search Term Discovery, overview of 10

## I

info\_gain score for terms 24, 44  
INPUT\_TERM\_PROP\_NAME pass-through  
definition 22  
filtering pre-tagged records 41, 42

## J

JSP code for rendering clusters 37

## L

LANG pass-through  
definition 22  
LANG\_PROP\_NAME pass-through  
definition 22

## M

MAX\_INPUT\_RECORDS pass-through  
definition 21  
minimal configuration for term extraction 15

## N

nav\_clusters.jsp sample file 49  
nav\_controls.jsp sample file 45  
noun phrases, size of 10

**O**

OUTPUT\_PROP\_NAME pass-through  
 definition 17  
 overlap properties for clusters 40

**P**

PARTIAL mode definition 21  
 partial updates for Term Discovery  
 mode setting 21  
 pass-throughs for term extraction  
 ALL\_TERMS\_OUTPUT\_PROP\_NAME 17  
 CORPUS\_MAX\_COVERAGE 24  
 CORPUS\_MIN\_COVERAGE 24  
 CORPUS\_MIN\_INFO\_GAIN 24  
 CORPUS\_MIN\_RECS 23  
 CORPUS\_REGEX\_KEEP 24  
 CORPUS\_REGEX\_SKIP 24  
 INPUT\_TERM\_PROP\_NAME 22  
 LANG\_PROP\_NAME 22  
 MAX\_INPUT\_RECORDS 21  
 minimal configuration 15  
 OUTPUT\_PROP\_NAME 17  
 RECORD\_FRACT\_OF\_MEDIAN 26  
 RECORD\_NTERMS 26  
 RECORD\_SPEC\_PROP\_NAME 16  
 TEXT\_PROP\_NAME 15  
 UPDATE\_MODE 20  
 pipeline for Term Discovery  
 exclude list record adapter 19  
 pre-tagged records for Term Discovery  
 corpus-level and record-level filtering 41  
 filtering for new terms only 42  
 filtering for pre-tagged and extracted terms 42  
 use cases 41

**R**

record adapters  
 exclude list 19  
 RECORD\_FRACT\_OF\_MEDIAN pass-through  
 definition 26  
 recommended setting 28  
 RECORD\_NTERMS pass-through  
 definition 26  
 recommended setting 28  
 RECORD\_SPEC\_PROP\_NAME pass-through 16  
 record-level filtering  
 best practices 28  
 configuration parameters 25

record-level filtering (*continued*)  
 for new terms only 42  
 for pre-tagged and extracted terms 42  
 for pre-tagged terms 41  
 setting scoring threshold 26  
 refinements  
 displaying in a Term Discovery dimension 40  
 displaying in nav\_controls.jsp 37  
 setting in controller file 36  
 regular expressions for term extraction 24  
 Relationship Discovery, overview of 9  
 relevant terms, definition of 10  
 restricting input records for term extraction 21

**S**

scoring threshold for record-level filtering, setting 26  
 search property for all extracted terms 17  
 singlet terms, eliminating 23  
 source property for term extraction 15  
 STATEFUL mode  
 definition 21  
 for tuning filtering pass-throughs 43  
 STATELESS mode for baseline updates 20  
 strategies to limit terms on records 26  
 Supplement objects for clusters 37

**T**

Term Discovery application  
 exclude list record adapter 19  
 global constants for UI 35  
 Term Discovery dimension  
 displaying records 40  
 displaying refinements 40  
 term extraction  
 format of source data 29  
 minimal configuration 15  
 overview 10  
 relevant terms 10  
 terms tagged on records, limiting 26  
 TEXT\_PROP\_NAME pass-through  
 definition 15  
 tuning strategy for clusters 33

**U**

UI for Term Discovery application 35  
 UPDATE\_MODE pass-through  
 values 20