

Oracle[®] Text

Application Developer's Guide

Release 9.0.1

June 2001

Part No. A90122-01

ORACLE[®]

Oracle Text Application Developer's Guide, Release 9.0.1

Part No. A90122-01

Copyright © 1996, 2001, Oracle Corporation. All rights reserved.

Primary Author: Colin McGregor

Contributors: Omar Alonso, Shamim Alpha, Steve Buxton, Chung-Ho Chen, Yun Cheng, Michele Cyran, Paul Dixon, Mohammad Faisal, Elena Huang, Garret Kaminaga, Ji Sun Kang, Bryn Llewellyn, Wesley Lin, Yasuhiro Matsuda, Gerda Shank, and Steve Yang.

The Programs (which include both the software and documentation) contain proprietary information of Oracle Corporation; they are provided under a license agreement containing restrictions on use and disclosure and are also protected by copyright, patent, and other intellectual and industrial property laws. Reverse engineering, disassembly, or decompilation of the Programs is prohibited.

The information contained in this document is subject to change without notice. If you find any problems in the documentation, please report them to us in writing. Oracle Corporation does not warrant that this document is error free. Except as may be expressly permitted in your license agreement for these Programs, no part of these Programs may be reproduced or transmitted in any form or by any means, electronic or mechanical, for any purpose, without the express written permission of Oracle Corporation.

If the Programs are delivered to the U.S. Government or anyone licensing or using the programs on behalf of the U.S. Government, the following notice is applicable:

Restricted Rights Notice Programs delivered subject to the DOD FAR Supplement are "commercial computer software" and use, duplication, and disclosure of the Programs, including documentation, shall be subject to the licensing restrictions set forth in the applicable Oracle license agreement. Otherwise, Programs delivered subject to the Federal Acquisition Regulations are "restricted computer software" and use, duplication, and disclosure of the Programs shall be subject to the restrictions in FAR 52.227-19, Commercial Computer Software - Restricted Rights (June, 1987). Oracle Corporation, 500 Oracle Parkway, Redwood City, CA 94065.

The Programs are not intended for use in any nuclear, aviation, mass transit, medical, or other inherently dangerous applications. It shall be the licensee's responsibility to take all appropriate fail-safe, backup, redundancy, and other measures to ensure the safe use of such applications if the Programs are used for such purposes, and Oracle Corporation disclaims liability for any damages caused by such use of the Programs.

Oracle is a registered trademark, and ConText, Oracle Text, Oracle8, Oracle8i, Oracle9i, Oracle Call Interface, PL/SQL, SQL*Plus, and SQL*Loader are trademarks or registered trademarks of Oracle Corporation. Other names may be trademarks of their respective owners.

Contents

Send Us Your Comments	ix
Preface.....	xi
1 Introduction to Oracle Text	
What is Oracle Text?	1-2
Types of Query Applications.....	1-2
Supported Document Formats	1-3
Theme Capabilities.....	1-3
Query Language and Operators.....	1-3
Document Services and Using a Thesaurus	1-4
Prerequisites For Building Your Query Application.....	1-4
Loading Your Text Table	1-5
Storing Text in the Text Table.....	1-6
Storing File Path Names	1-6
Storing URLs	1-6
Storing Associated Document Information	1-7
Supported Column Types	1-7
Supported Document Formats	1-7
Loading Methods.....	1-8
Indexing Your Documents	1-9
Type of Index	1-9
Creating a CONTEXT Index	1-10
Creating a CTXCAT Index	1-12
Creating a CTXRULE Index.....	1-12

Index Maintenance	1-12
A Simple Text Query Application	1-14
Querying your Index	1-16
Querying with CONTAINS.....	1-16
Structured Field Searching	1-17
Thesaural Queries.....	1-18
Document Section Searching.....	1-18
Other Query Features.....	1-18
Presenting the Hitlist	1-20
Hitlist Example.....	1-20
Presenting Structured Fields.....	1-22
Ordering the Hitlist	1-22
Presenting Document Hit Count	1-22
Document Presentation and Highlighting.....	1-23
Highlighting Example.....	1-24
Document List of Themes Example	1-25
Gist Example.....	1-26

2 Indexing

About Oracle Text Indexes	2-2
Structure of the Oracle Text CONTEXT Index.....	2-2
The Oracle Text Indexing Process	2-4
Partitioned Tables and Indexes.....	2-6
Parallel Indexing	2-6
Limitations for Indexing	2-7
Considerations For Indexing	2-8
Type of Index.....	2-9
Location of Text.....	2-10
Document Formats and Filtering	2-11
Bypassing Rows for Indexing	2-12
Document Character Set	2-12
Document Language	2-13
Indexing Special Characters.....	2-13
Case-Sensitive Indexing and Querying.....	2-15
Language Specific Features	2-15

Fuzzy Matching and Stemming	2-17
Better Wildcard Query Performance	2-17
Document Section Searching	2-18
Stopwords and Stopthemes	2-18
Index Creation	2-20
Procedure for Creating a CONTEXT Index	2-20
Creating Preferences	2-21
Creating Section Groups for Section Searching	2-25
Using Stopwords and Stoplists.....	2-25
Creating an Index	2-27
Creating a CONTEXT Index	2-27
Creating a CTXCAT Index	2-29
Creating a CTXRULE Index.....	2-32
Index Maintenance	2-34
Viewing Index Errors.....	2-34
Dropping an Index	2-34
Resuming Failed Index	2-34
Rebuilding an Index.....	2-35
Dropping a Preference	2-35
Managing DML Operations for a CONTEXT Index	2-36
Viewing Pending DML.....	2-36
Synchronizing the Index.....	2-36
Index Optimization	2-37

3 Querying

Overview of Queries	3-2
Querying with CONTAINS	3-2
Querying with CATSEARCH	3-4
Querying with MATCHES.....	3-5
Word and Phrase Queries	3-7
ABOUT Queries and Themes	3-8
Query Expressions.....	3-9
Case-Sensitive Searching	3-10
Query Feedback	3-11
Query Explain Plan	3-11

Query Operators for CONTAINS	3-13
ABOUT Query	3-13
Logical Operators.....	3-13
Section Searching	3-15
Proximity Queries with NEAR Operator	3-15
Fuzzy, Stem, Soundex, Wildcard and Thesaurus Expansion Operators.....	3-15
Stored Query Expressions	3-15
Calling PL/SQL Functions in CONTAINS.....	3-16
Query Operators for CATSEARCH	3-18
Optimizing for Response Time	3-19
Retrieving a Range of Documents.....	3-19
Counting Hits	3-21
SQL Count Hits Example.....	3-21
Counting Hits with a Structured Predicate.....	3-21
PL/SQL Count Hits Example	3-21

4 Document Presentation

Highlighting Query Terms	4-2
Text highlighting.....	4-2
Theme Highlighting	4-2
CTX_DOC Highlighting Procedures	4-2
Obtaining List of Themes, Gists, and Theme Summaries	4-4
List of Themes	4-4
Gist and Theme Summary.....	4-5

5 Query Tuning

Optimizing Queries with Statistics	5-2
Collecting Statistics.....	5-2
Re-Collecting Statistics.....	5-3
Deleting Statistics.....	5-4
Optimizing Queries for Response Time	5-5
Better Response Time with FIRST_ROWS	5-5
Better Response Time with CHOOSE	5-6
Optimizing Queries for Throughput	5-8
CHOOSE and ALL ROWS Modes.....	5-8

	FIRST_ROWS Mode	5-8
	Tuning Queries with Blocking Operations	5-9
6	Document Section Searching	
	About Document Section Searching.....	6-2
	Enabling Section Searching	6-2
	Section Types.....	6-5
	HTML Section Searching	6-10
	Creating HTML Sections	6-10
	Searching HTML Meta Tags	6-10
	XML Section Searching	6-12
	Automatic Sectioning.....	6-12
	Attribute Searching	6-12
	Creating Document Type Sensitive Sections.....	6-13
	Path Section Searching.....	6-14
7	Working With a Thesaurus	
	Overview of Thesauri	7-2
	Thesaurus Creation and Maintenance.....	7-2
	Case-sensitive Thesauri	7-3
	Case-insensitive Thesauri.....	7-3
	Default Thesaurus	7-4
	Supplied Thesaurus.....	7-4
	Defining Thesaural Terms.....	7-6
	Defining Synonyms.....	7-6
	Defining Hierarchical Relations	7-6
	Using a Thesaurus in a Query Application	7-8
	Loading a Custom Thesaurus and Issuing Thesaural Queries.....	7-8
	Augmenting Knowledge Base with Custom Thesaurus	7-9
	About the Supplied Knowledge Base.....	7-12
	Adding a Language-Specific Knowledge Base	7-13
8	Administration	
	Oracle Text Users and Roles	8-2

CTXSYS User	8-2
CTXAPP Role	8-2
Granting Roles and Privileges to Users.....	8-2
DML Queue.....	8-3
The CTX_OUTPUT Package.....	8-4
Servers.....	8-5
Administration Tool	8-6

A CONTEXT Query Application

Web Query Application Overview	A-2
The PSP Web Application	A-2
Web Application Prerequisites	A-3
Building the Web Application	A-3
Web Application Sample Code	A-6
loader.ctl.....	A-6
loader.dat	A-6
search_htmlservices.sql.....	A-7
search_html.psp	A-9

Index

Send Us Your Comments

Oracle Text Application Developer's Guide, Release 9.0.1

Part No. A90122-01

Oracle Corporation welcomes your comments and suggestions on the quality and usefulness of this document. Your input is an important part of the information used for revision.

- Did you find any errors?
- Is the information clearly presented?
- Do you need more information? If so, where?
- Are the examples correct? Do you need more examples?
- What features did you like most?

If you find any errors or have any other suggestions for improvement, please indicate the document title and part number, and the chapter, section, and page number (if available). You can send comments to us in the following ways:

- Electronic mail: infodev_us@us.oracle.com
- FAX: (650) 506-7227 Attn: Server Technologies Documentation Manager
- Postal service:

Oracle Corporation
Server Technologies Documentation
500 Oracle Parkway, Mailstop 4op11
Redwood Shores, CA 94065
USA

If you would like a reply, please give your name, address, telephone number, and (optionally) electronic mail address.

If you have problems with the software, please contact your local Oracle Support Services.

Preface

This guide explains how to build query applications with Oracle Text. This preface contains these topics:

- [Audience](#)
- [Organization](#)
- [Related Documentation](#)
- [Conventions](#)
- [Documentation Accessibility](#)

Audience

Oracle Text Application Developer's Guide is intended for users who perform the following tasks:

- Develop Oracle Text applications.
- Administer Oracle Text installations.

To use this document, you need to have experience with the Oracle object relational database management system, SQL, SQL*Plus, and PL/SQL.

Organization

This document contains:

Chapter 1, "Introduction to Oracle Text"

This chapter introduces the basic features of Oracle Text. It also explains how to build a basic query application using Oracle Text.

Chapter 2, "Indexing"

This chapter describes how to index your document set. It discusses considerations for indexing as well as how to create CONTEXT, CTXCAT, and CTXRULE indexes.

Chapter 3, "Querying"

This chapter describes how to query your document set. It gives examples for using the CONTAINS, CATSEARCH, and MATCHES operators.

Chapter 4, "Document Presentation"

This chapter describes how to present documents to the user of your query application.

Chapter 5, "Query Tuning"

This chapter describes how to tune your queries to improve response time and throughput.

Chapter 6, "Document Section Searching"

This chapter describes how to enable section searching in HTML and XML.

Chapter 7, "Working With a Thesaurus"

This chapter describes how to work with a thesaurus in your application. It also describes how to augment your knowledge with a thesaurus.

Chapter 8, "Administration"

This chapter describes Oracle Text administration.

Appendix A, "CONTEXT Query Application"

This chapter describes an Oracle Text example web application.

Related Documentation

For more information about *Oracle Text*, see:

- *Oracle Text Reference*

For more information about Oracle9i, see:

- *Oracle9i Database Concepts*
- *Oracle9i Database Administrator's Guide*
- *Oracle9i Database Utilities*
- *Oracle9i Database Performance Guide and Reference*
- *Oracle9i SQL Reference*
- *Oracle9i Database Reference*
- *Oracle9i Application Developer's Guide - Fundamentals*
- *Oracle9i Application Developer's Guide - XML*

For more information about PL/SQL, see:

- *PL/SQL User's Guide and Reference*

In North America, printed documentation is available for sale in the Oracle Store at

<http://oraclestore.oracle.com/>

Customers in Europe, the Middle East, and Africa (EMEA) can purchase documentation from

<http://www.oraclebookshop.com/>

Other customers can contact their Oracle representative to purchase printed documentation.

To download free release notes, installation documentation, code samples, white papers, or other collateral, please visit the Oracle Technology Network (OTN). You must register online before using OTN; registration is free and can be done at

<http://technet.oracle.com/membership/index.htm>

If you already have a username and password for OTN, then you can go directly to the documentation section of the OTN Web site at

<http://technet.oracle.com/docs/index.htm>

You can obtain Oracle Text technical information, collateral, code samples, training slides and other material at:

<http://technet.oracle.com/products/text>

Conventions

This section describes the conventions used in the text and code examples of the this documentation set. It describes:

- [Conventions in Text](#)
- [Conventions in Code Examples](#)

Conventions in Text

We use various conventions in text to help you more quickly identify special terms. The following table describes those conventions and provides examples of their use.

Convention	Meaning	Example
Bold	Bold typeface indicates terms that are defined in the text or terms that appear in a glossary, or both.	The C datatypes such as ub4 , sword , or OCINumber are valid. When you specify this clause, you create an index-organized table .
<i>Italics</i>	Italic typeface indicates query terms, book titles, emphasis, syntax clauses, or placeholders.	<i>Oracle9i Database Concepts</i> You can specify the <i>parallel_clause</i> . Run <i>Uold_release</i> .SQL where <i>old_release</i> refers to the release you installed prior to upgrading.

Convention	Meaning	Example
UPPERCASE monospace (fixed-width font)	Uppercase monospace typeface indicates elements supplied by the system. Such elements include parameters, privileges, datatypes, RMAN keywords, SQL keywords, SQL*Plus or utility commands, packages and methods, as well as system-supplied column names, database objects and structures, user names, and roles.	<p>You can specify this clause only for a NUMBER column.</p> <p>You can back up the database using the BACKUP command.</p> <p>Query the TABLE_NAME column in the USER_TABLES table in the data dictionary view.</p> <p>Specify the ROLLBACK_SEGMENTS parameter.</p> <p>Use the DBMS_STATS.GENERATE_STATS procedure.</p>
lowercase monospace (fixed-width font)	Lowercase monospace typeface indicates executables and sample user-supplied elements. Such elements include computer and database names, net service names, and connect identifiers, as well as user-supplied database objects and structures, column names, packages and classes, user names and roles, program units, and parameter values.	<p>Enter sqlplus to open SQL*Plus.</p> <p>The department_id, department_name, and location_id columns are in the hr.departments table.</p> <p>Set the QUERY_REWRITE_ENABLED initialization parameter to true.</p> <p>Connect as oe user.</p>

Conventions in Code Examples

Code examples illustrate SQL, PL/SQL, SQL*Plus, or other command-line statements. They are displayed in a monospace (fixed-width) font and separated from normal text as shown in this example:

```
SELECT username FROM dba_users WHERE username = 'MIGRATE';
```

The following table describes typographic conventions used in code examples and provides examples of their use.

Convention	Meaning	Example
[]	Brackets enclose one or more optional items. Do not enter the brackets.	DECIMAL (digits [, precision])
{ }	Braces enclose two or more items, one of which is required. Do not enter the braces.	{ENABLE DISABLE}
	A vertical bar represents a choice of two or more options within brackets or braces. Enter one of the options. Do not enter the vertical bar.	{ENABLE DISABLE} [COMPRESS NOCOMPRESS]

Convention	Meaning	Example
...	Horizontal ellipsis points indicate either: <ul style="list-style-type: none"> That we have omitted parts of the code that are not directly related to the example That you can repeat a portion of the code 	<pre>CREATE TABLE ... AS subquery; SELECT col1, col2, ... , coln FROM employees;</pre>
.	Vertical ellipsis points indicate that we have omitted several lines of code not directly related to the example.	
Other notation	You must enter symbols other than brackets, braces, vertical bars, and ellipsis points as it is shown.	<pre>acctbal NUMBER(11,2); acct CONSTANT NUMBER(4) := 3;</pre>
<i>Italics</i>	Italicized text indicates variables for which you must supply particular values.	<pre>CONNECT SYSTEM/<i>system_password</i></pre>
UPPERCASE	Uppercase typeface indicates elements supplied by the system. We show these terms in uppercase in order to distinguish them from terms you define. Unless terms appear in brackets, enter them in the order and with the spelling shown. However, because these terms are not case sensitive, you can enter them in lowercase.	<pre>SELECT last_name, employee_id FROM employees; SELECT * FROM USER_TABLES; DROP TABLE hr.employees;</pre>
lowercase	Lowercase typeface indicates programmatic elements that you supply. For example, lowercase indicates names of tables, columns, or files.	<pre>SELECT last_name, employee_id FROM employees; sqlplus hr/hr</pre>

Documentation Accessibility

Oracle's goal is to make our products, services, and supporting documentation accessible to the disabled community with good usability. To that end, our documentation includes features that make information available to users of assistive technology. This documentation is available in HTML format, and contains markup to facilitate access by the disabled community. Standards will continue to evolve over time, and Oracle is actively engaged with other market-leading technology vendors to address technical obstacles so that our documentation can be accessible to all of our customers. For additional information, visit the Oracle Accessibility Program Web site at

<http://www.oracle.com/accessibility/>

JAWS, a Windows screen reader, may not always correctly read the code examples in this document. The conventions for writing code require that closing braces should appear on an otherwise empty line; however, JAWS may not always read a line of text that consists solely of a bracket or brace.

Introduction to Oracle Text

This chapter introduces the main features of Oracle Text. It is provided to help you get started with indexing, querying, and document presentation.

The following topics are covered:

- [What is Oracle Text?](#)
- [Loading Your Text Table](#)
- [Indexing Your Documents](#)
- [A Simple Text Query Application](#)
- [Querying your Index](#)
- [Presenting the Hitlist](#)
- [Document Presentation and Highlighting](#)

What is Oracle Text?

Oracle Text is a tool that enables you to build text query applications and document classification applications. Oracle Text provides indexing, word and theme searching, and viewing capabilities for text.

Types of Query Applications

You can build two types of applications with Oracle Text:

- Text Query Application
- Document Classification Application

Text Query Applications

The purpose of a text query application is to enable users to find text that contains one or more search terms. The text is usually a collection of documents. A good application can index and search common document formats such as HTML, XML, plain text, or Microsoft Word. For example, an application with a browser interface might enable users to query a company website consisting of HTML files, returning those files that match a query.

To build a text query application, you can create either a `context` or `ctxcat` index and query the index with `CONTAINS` or `CATSEARCH` respectively.

Document Classification Applications

A document classification application is one that classifies an incoming stream of documents based on its content. They are also known as document routing or filtering applications. For example, an online news agency might need to classify its incoming stream of articles as they arrive into categories such as politics, crime, and sports.

Oracle Text enables you to build these applications with the `CTXRULE` index type. This index type indexes the rules (queries) that define each class. When documents arrive, the `MATCHES` operator can be used to match each document with the rules that select it.

Note: Oracle Text supports document classification for only plain text, XML, and HTML documents.

See Also: ["Indexing Your Documents"](#) in this chapter for more information about these index types.

Supported Document Formats

For text query applications, Oracle Text supports most document formats for indexing and querying, including plain text, HTML and formatted documents such as Microsoft Word.

For document classification application, Oracle Text supports classifying plain text, HTML, and XML documents.

Theme Capabilities

With Oracle Text, you can search on document themes if your language is English and French. To do so, you use the ABOUT operator. For example, you can search for all documents that are about the concept *politics*. Documents returned might be about elections, governments, or foreign policy. The documents need not contain the word *politics* to score hits.

Theme information is derived from the supplied knowledge base, which is a hierarchical listing of categories and concepts. As the supplied knowledge base is a general view of the world, you can add to it new industry-specific concepts. With an augmented knowledge base, the system can process document themes more intelligently and so improve the accuracy of your theme searching.

With the supplied PL/SQL packages, you can also obtain document themes programmatically.

See Also: *Oracle Text Reference* to learn more about the ABOUT operator.

Themes in Other Languages

You can enable theme capabilities such as ABOUT queries in other languages besides English and French by loading a language-specific knowledge base.

See Also: [Adding a Language-Specific Knowledge Base](#) in [Chapter 7, "Working With a Thesaurus"](#).

Query Language and Operators

To query, you use the SQL SELECT statement. Depending on your index, you can query text with either the CONTAINS operator, which is used with the context

index, or the CATSEARCH operator, which is used with the ctxcat index. You use these operators in the WHERE clause of the SELECT statement as follows:

```
SELECT SCORE(1) title from news WHERE CONTAINS(text, 'oracle', 1) > 0;
```

To classify single documents, you use the MATCHES operator with a ctxrule index.

For text querying with the CONTAINS operator, Oracle Text provides a rich query language with operators that enable you to issue variety of queries including simple word queries, ABOUT queries, logical queries, wildcard and thesaural expansion queries.

The CATSEARCH operator also supports some of the operations available with CONTAINS.

See Also: [Chapter 3, "Querying"](#)

Document Services and Using a Thesaurus

You can also use the supplied Oracle Text PL/SQL packages for advanced features such as document presentation and thesaurus maintenance.

See Also:

[Chapter 7, "Working With a Thesaurus"](#)

[Chapter 4, "Document Presentation"](#)

Prerequisites For Building Your Query Application

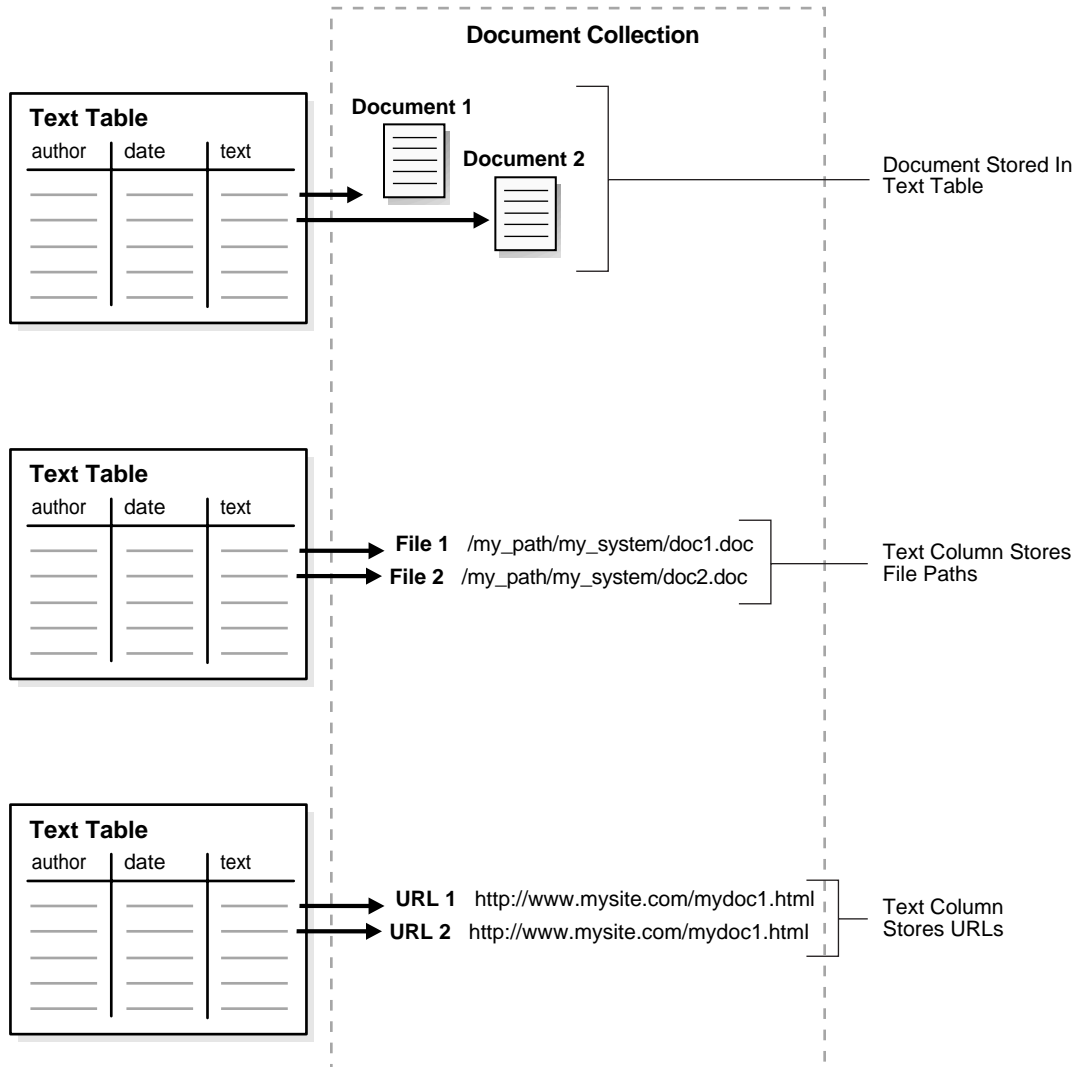
To build an Oracle Text query application, you must have the following:

- a populated text table
- an Oracle Text index

The following sections describe these prerequisites and also describe the main features of a generic text query application.

Loading Your Text Table

Figure 1-1



The basic prerequisite for an Oracle Text query application is to have a populated text table. The text table is where you store information about your document collection and is required for indexing.

You can populate rows in your text table with one of the following elements:

- text information (can be documents or text fragments)
- path names of documents in your file system
- URLs that specify world-wide web documents

Figure 1-1 illustrates these different methods.

By default, the indexing operation expects your document text to be directly loaded in your text table, which is the first method above.

However, you can specify the other ways of identifying your documents such as with filenames or with URLs using the corresponding data storage indexing preference.

Storing Text in the Text Table

You can store documents in your text table in different ways.

You can store documents in one column using the `DIRECT_DATASTORE` data storage type or over a number of columns using the `MULTI_COLUMN_DATASTORE` type. When your text is stored over a number of columns, Oracle concatenates the columns into a virtual document for indexing.

You can also create master-detail relationships for your documents, where one document can be stored across a number of rows. To create master-detail index, use the `DETAIL_DATASTORE` data storage type.

In your text table, you can also store short text fragments such as names, descriptions, and addresses over a number of columns to create a catalog index.

You can also store your text in a nested table using the `NESTED_DATASTORE` type.

Oracle Text supports the indexing of the XMLType which you can use to store XML documents.

Storing File Path Names

In your text table, you can store path names to files stored in your file system. When you do so, use the `FILE_DATASTORE` preference type during indexing.

Storing URLs

You can store URL names to index web-sites. When you do so, use the `URL_DATASTORE` preference type during indexing.

Storing Associated Document Information

In your text table, you can create additional columns to store structured information that your query application might need, such as primary key, date, description, or author.

Format and Character Set Columns

If your documents are of mixed formats or of mixed character sets, you can create the following additional columns:

- Format column to record format (TEXT or BINARY) to help filtering during indexing. You can also use to format column to ignore rows for indexing by setting the format column to IGNORE. This is useful for bypassing rows that contain data incompatible with text indexing such as images.
- Character set column to record document character set on a per row basis.

When you create your index, you must specify the name of the format or character set column in the parameter clause of CREATE INDEX.

Supported Column Types

With Oracle Text, you can create a CONTEXT index with columns of type VARCHAR2, CLOB, BLOB, CHAR, BFILE, and XMLType.

Note: The column types NCLOB, DATE and NUMBER cannot be indexed.

Supported Document Formats

Because the system can index most document formats including HTML, PDF, Microsoft Word, and plain text, you can load any supported type into the text column.

When you have mixed formats in your text column, you can optionally include a format column to help filtering during indexing. With the format column you can specify whether a document is binary (formatted) or text (non-formatted such as HTML).

See Also: *Oracle Text Reference* for more information about the supported document formats.

Loading Methods

The following sections describe different methods of loading information into a text column.

INSERT Statement

You can use the SQL INSERT statement to load text to a table.

The following example creates a table with two columns, `id` and `text`, using CREATE TABLE. The example populates the table with the INSERT statement. This example makes the `id` column the primary key, which is the required constraint for a Text table. The `text` column is VARCHAR2:

```
CREATE TABLE docs (id number primary key, text varchar2(80));
```

To populate this table, use the INSERT statement as follows:

```
INSERT into docs values(1, 'this is the text of the first document');  
INSERT into docs values(12, 'this is the text of the second document');
```

Loading Text from File-System

In addition to the INSERT statement, Oracle enables you to load text data (this includes documents, pointers to documents, and URLs) into a table from your file-system using other automated methods, including

- SQL*Loader
- DBMS_LOB.LOADFROMFILE() PL/SQL procedure to load LOBs from BFILES
- Oracle Call Interface

See Also:

[Appendix A, "CONTEXT Query Application"](#) for a SQL*Loader example.

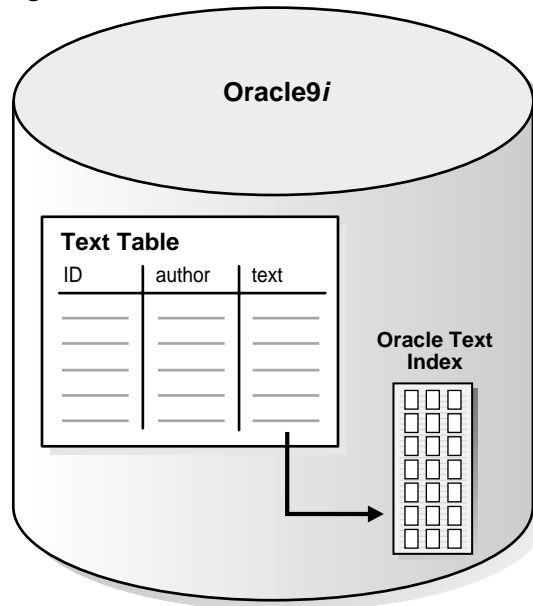
Oracle9i Supplied PL/SQL Packages Reference For more information about the DBMS_LOB package.

Oracle9i Application Developer's Guide - Large Objects (LOBs) for more information about working with LOBs.

Oracle Call Interface Programmer's Guide for more information about Oracle Call Interface.

Indexing Your Documents

Figure 1–2



To query your document collection, you must first index the text column of your text table. Indexing breaks your text into tokens, which are usually words. This creates a CONTEXT index, which records each token and the documents that contain it. An inverted index as such allows for querying on words and phrases. [Figure 1–2](#) shows a text table within Oracle9i and its associated Oracle Text index.

Type of Index

Oracle Text supports the creation of three types of indexes depending on your application and text source. You use the CREATE INDEX statement to create all Oracle Text index types.

The following table describes these indexes and the type of applications you can build with them. The third column shows which query operator to use with the index.

Index Type	Application Type	Query Operator
CONTEXT	Use this index to build a text retrieval application when your text consists of large coherent documents. You can index documents of different formats such as MS Word, HTML, XML, or plain text. With a context index, you can customize your index in a variety of ways.	CONTAINS
CTXCAT	Use this index type to improve mixed query performance. Suitable for querying small text fragments with structured criteria like dates, item names, and prices that are stored across columns.	CATSEARCH
CTXRULE	Use a CTXRULE index to build a document classification application. The CTXRULE index is an index created on a table of queries, where each query has a classification. Single documents (plain text, HTML, or XML) can be classified using the MATCHES operator.	MATCHES

Creating a CONTEXT Index

Once your text data is loaded in a table, you can use `CREATE INDEX` to create a `context` index. When you create an index and specify no parameter clause, an index is created with default parameters.

For example, the following command creates a `context` index called `myindex` on the `text` column in the `docs` table:

```
CREATE INDEX myindex ON docs(text) INDEXTYPE IS CTXSYS.CONTEXT;
```

General Defaults for All Languages

When you use `CREATE INDEX` to create a context index without explicitly specifying parameters, the system does the following for all languages by default:

- Assumes that the text to be indexed is stored directly in a text column. The text column can be of type `CLOB`, `BLOB`, `BFILE`, `VARCHAR2`, `XMLType`, or `CHAR`.

- Detects the column type and uses filtering for binary column types. Most document formats are supported for filtering. If your column is plain text, the system does not use filtering.

Note: For document filtering to work correctly in your system, you must ensure that your environment is set up correctly to support the Inso filter.

To learn more about configuring your environment to use the Inso filter, see *Oracle Text Reference*.

- Assumes the language of text to index is the language you specify in your database setup.
- Uses the default stoplist for the language you specify in your database setup. Stoplists identify the words that the system ignores during indexing.
- Enables fuzzy and stemming queries for your language, if this feature is available for your language.

You can always change the default indexing behavior by creating your own preferences and specifying these custom preferences in the parameter clause of CREATE INDEX.

Customizing Your CONTEXT Index

Using the parameter clause with CREATE INDEX, you can customize your context index. For example, in the parameter clause, you can specify where your text is stored, how you want it filtered for indexing, and whether sections should be created.

To index a set of HTML files loaded in the text column `htmlfile`, you can issue the CREATE INDEX statement, specifying `datastore`, `filter` and `section group` parameters as follows:

```
CREATE INDEX myindex ON doc(htmlfile) INDEXTYPE IS ctxsys.context PARAMETERS
('datastore ctxsys.default_datastore filter ctxsys.null_filter section group
ctxsys.html_section_group');
```

See Also: "[Considerations For Indexing](#)" in Chapter 2, "[Indexing](#)" for more information about the different ways you can create an index.

Oracle Text Reference for more information on the CREATE INDEX statement.

Creating a CTXCAT Index

A CTXCAT index is an index optimized for mixed queries. You can create this type of index when you store small documents or text fragments and associated structured information. To query this index, you use the CATSEARCH operator and specify a structured clause, if any. Query performance with a CTXCAT index is usually better for structured queries than with a CONTEXT index.

See Also: ["Creating a CTXCAT Index" in Chapter 2, "Indexing"](#) for a complete example.

Creating a CTXRULE Index

You create a CTXRULE index to build a document classification application in which an incoming stream of documents is routed according content. You define the classification rules as queries which you index. You use the MATCHES operator to classify single documents.

See Also: ["Creating a CTXRULE Index" in Chapter 2, "Indexing"](#) for a complete example.

Index Maintenance

Index maintenance is necessary after your application inserts, updates, or deletes documents in your base table.

If your base table is static, that is, you do no updating, inserting or deleting of documents after your initial index, you do not need to maintain your index.

However, if you perform DML operations (inserts, updates, or deletes) on your base table, you must update your index. You can synchronize your index manually with CTX_DDL.SYNC_INDEX.

The following example synchronizes the index `myindex` with 2 megabytes of memory:

```
begin
  ctx_ddl.sync_index('myindex', '2M');
end;
```

If you synchronize your index regularly, you might also consider optimizing your index to reduce fragmentation and to remove old data.

See Also: ["Managing DML Operations for a CONTEXT Index"](#) in [Chapter 2, "Indexing"](#) for more information about synchronizing and optimizing the index.

A Simple Text Query Application

A typical search application allows the user to enter a query. The application executes the query and returns a list of documents, usually ranked by relevance, that satisfy the query. The application enables the user to view one or more documents in the returned hitlist.

For example, a application portal might index URLs (HTML files) on the worldwide web and provide query capabilities across the set of indexed URLs. Query hitlists are composed of URLs that the user can visit.

Figure 1–3 Flowchart of a Typical Query Application

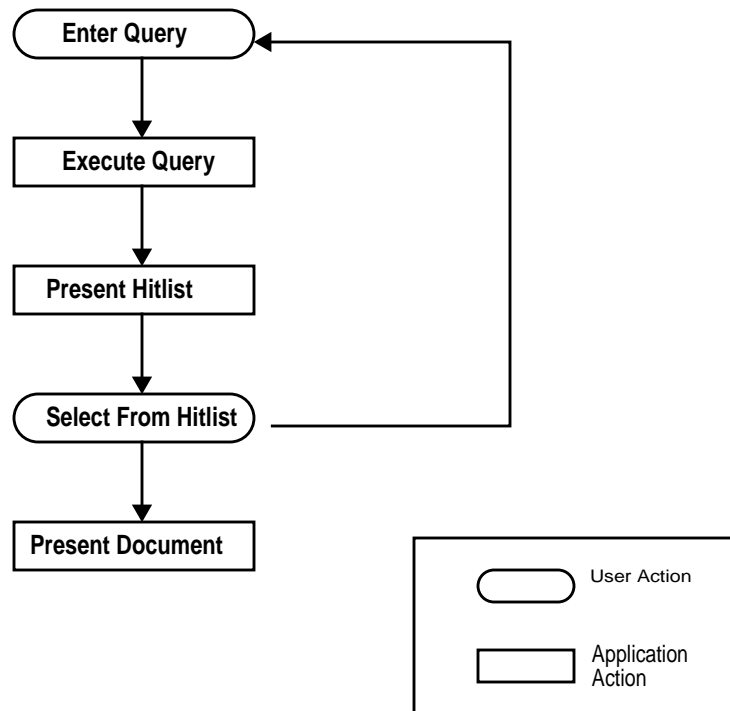


Figure 1–3 illustrates the flowchart of how a user interacts with a simple query application. The figure shows the steps required to enter the query and to view the results. Rectangular boxes indicate application tasks and oval boxes indicate user-tasks.

As shown, the a query application can be modeled according to the following steps:

1. user enters query
2. application executes query
3. application presents hitlist
4. user selects document from hitlist
5. application presents document to user for viewing

The rest of this chapter explains how you can accomplish these steps with Oracle Text.

See Also: [Appendix A, "CONTEXT Query Application"](#) for a description and of a simple web query application.

Querying your Index

With Oracle Text, you use the CONTAINS operator to query a `context` index. This is the most common operator and index used to build query applications.

For more advanced applications, you use the CATSEARCH operator to query a `ctxcat` index, and you use the MATCHES operator to query the `ctxrule` index.

Querying with CONTAINS

You can use CONTAINS to retrieve documents that contain a word or phrase. Your document must be indexed before you can issue a CONTAINS query.

Use the CONTAINS operator in a SELECT statement. With CONTAINS, you can issue two types of queries:

- word query
- ABOUT query

You can also optimize queries for better response time to obtain the top n hits. The following sections give an overview of these query scenarios.

Word Query Example

A word query is a query on the exact word or phrase you enter between the single quotes in the CONTAINS or CATSEARCH operator.

The following example finds all the documents in the `text` column that contain the word *oracle*. The score for each row is selected with the SCORE operator using a label of 1:

```
SELECT SCORE(1) title FROM news WHERE CONTAINS(text, 'oracle', 1) > 0;
```

In your query expression, you can use text operators such as AND and OR to achieve different results. You can also add structured predicates to the WHERE clause.

See Also: *Oracle Text Reference* for more information about the different operators you can use in queries.

You can count the hits to a query using the SQL COUNT(*) statement, or CTX_QUERY.COUNT_HITS.

ABOUT Query Example

In all languages, ABOUT queries increases the number of relevant documents returned by a query.

In English and French, ABOUT queries can use the theme component of the index, which is created by default. As such, this operator returns documents based on the concepts of your query, not only the exact word or phrase you specify.

For example, the following query finds all the documents in the *text* column that are about the subject *politics*, not just the documents that contain the word *politics*:

```
SELECT SCORE(1) title FROM news WHERE CONTAINS(text, 'about(politics)', 1) > 0;
```

See Also: *Oracle Text Reference* to learn more about the ABOUT operator.

Optimizing Query for Response Time

You can optimize any CONTAINS query (word or ABOUT) for response time in order to retrieve the highest ranking hits in a result set in the shortest time possible. Optimizing for response time is useful in a web-based search application.

See Also: ["Optimizing Queries for Response Time"](#) in [Chapter 5](#), ["Query Tuning"](#).

Structured Field Searching

Your application interface can give the user the option of querying on structured fields related to the text such as item description, author, or date as a means of further limiting the search criteria.

You can issue structured searches with CONTAINS using a structured clause in the SELECT statement. However, for optimal performance, consider creating a *ctxcat* index which gives better performance for structured queries with the CATSEARCH operator.

Your application can also present the structured information related to each document in the hitlist.

See Also: ["Creating a CTXCAT Index"](#) in [Chapter 2](#), ["Indexing"](#) for more information about creating a *ctxcat* index to improve structured queries with CATSEARCH.

Thesaural Queries

Oracle Text enables you to define a thesaurus for your query application.

Defining a custom thesaurus allows you to process queries more intelligently. Since users of your application might not know which words represent a topic, you can define synonyms or narrower terms for likely query terms. You can use the thesaurus operators to expand your query with thesaurus terms.

See Also: [Chapter 7, "Working With a Thesaurus"](#)

Document Section Searching

Section searching enables you to narrow text queries down to sections within documents.

Section searching can be implemented when your documents have internal structure, such as HTML and XML documents. For example, you can define a section for the <H1> tag that allows you to query within this section using the WITHIN operator.

You can set the system to automatically create sections from XML documents.

You can also define attribute sections to search attribute text in XML documents.

Note: Section searching is supported for only word queries with a `context` index type.

See Also: [Chapter 6, "Document Section Searching"](#)

Other Query Features

In your query application, you can use other query features such as proximity searching. The following table lists some of these features.

Feature	Description	Implement With
Proximity Searching	Enables searches for words near one another.	NEAR operator at query-time.
Stemming	Enables searching for words with same root as specified term.	\$ operator at query-time.

Feature	Description	Implement With
Fuzzy Searching	Search for words that have similar spelling to specified term.	fuzzy operator at query-time.
Case Sensitive Searching	Enables case-sensitive searches.	Enable with basic_lexer at index-time.
Base Letter Conversion	Queries match words with or without diacritical marks such as tildes, accents, and umlauts. For example, with a Spanish base-letter index, a query of <i>energía</i> matches documents containing <i>energía</i> and <i>energía</i> .	Enable with basic_lexer at index-time.
Word Decomposing (German and Dutch)	Enables searching on words that contain specified term as sub-composite.	Enable with basic_lexer at index-time.
Alternate Spelling (German, Dutch, and Swedish)	Enables searches on alternate spellings of words.	Enable with basic_lexer at index-time.
Query Explain Plan	Generate query parse information.	CTX_QUERY.EXPLAIN
Hierarchical Query Feedback	Generate broader term, narrower term and related term information for a query.	CTX_QUERY.HFEEDBACK
Browse index	Browse the words around a seed word in the index.	CTX_QUERY.BROWSE_WORDS
Count hits	Count the number of hits in a query	CTX_QUERY.COUNT_HITS
Stored Query Expression	Stores a query expression	CTX_QUERY.STORE_SQE
Thesaural Queries	Use a thesaurus to expand queries.	Thesaurus operators such as SYN and BT as well as the ABOUT operator. Use CTX_THES package to maintain thesaurus.

Presenting the Hitlist

After executing the query, query applications typically present a hitlist of all documents that satisfy the query along with a relevance score. This list can be a list of document titles or URLs depending on your document set.


Your application presents a hitlist in one or more of the following ways:

- show documents ordered by score
- show structured fields related to document, such as title or author
- show document hit count

Hitlist Example

[Figure 1-4](#) is a screen shot of a query application presenting the hitlist to the user.















Figure 1-4 Query Application Presenting Hitlist



Search for

Page 1 of 4 Oracle found 34 results for Java and XML

[\[Next\]](#)

- 1 [Oracle Corp. Document](#)  
Gist: Oracle has built new functionality on the object support of Oracle8 to produce O ...
Precision score: 100%, **Last modified:** 13-MAR-00, **Page size:** 56731 bytes
Keywords: Oracle8i, paradigms, Internet, CYBERSMARTS, delivery, technology, systems, Oracl ...
- 2 [White Paper](#)  
Gist: Supporting standard EJB and CORBA deployment architectures, Oracle Business Comp ...
Precision score: 67%, **Last modified:** 13-MAR-00, **Page size:** 29053 bytes
Keywords: Oracle Corporation, Java, Oracle Business Components, logic, relation, writing, ...
- 3 [Oracle JDeveloper Suite - Technical Documentation - Partners](#)  
Gist: Partner Software Solutions Halcyon Software produces a suite of products for c ...
Precision score: 43%, **Last modified:** 13-MAR-00, **Page size:** 6464 bytes
Keywords: partners, Oracle JDeveloper, computer software, Java, Oracle Corporation, conver ...
- 4 [Oracle JDeveloper - Technical Resources - Documentation](#)  
Gist: Using the Oracle Business Components framework, the developer can focus on writi ...
Precision score: 19%, **Last modified:** 13-MAR-00, **Page size:** 7254 bytes
Keywords: Java, Oracle Corporation, Oracle JDeveloper, technology, logic, Oracle Business ...
- 5 [Oracle Corp. Document](#)  
Gist: Oracle8i Adds InterMedia, iFS, and XML Support to Provide Leading Platform for I ...
Precision score: 18%, **Last modified:** 02-MAR-00, **Page size:** 10050 bytes
Keywords: Oracle Corporation, management, contents, Oracle8i, Internet, technology, Oracle ...
- 6 [Oracle Corp. Document](#)  
Gist: "Over the last decade corporations have built meta-data islands. IT departments ...
Precision score: 18%, **Last modified:** 02-MAR-00, **Page size:** 4711 bytes
Keywords: Oracle Corporation, Oracle Repository, management, COITS, construction, meta, fun ...
- 7 [Oracle Corp. Document](#)  
Gist: "Over the last decade corporations have built meta-data islands. IT departments ...
Precision score: 18%, **Last modified:** 02-MAR-00, **Page size:** 4711 bytes
Keywords: Oracle Corporation, Oracle Repository, management, COITS, construction, meta, fun ...

Presenting Structured Fields

Structured columns related to the text column can help identify documents. When you present the hitlist, you can show related columns such as document titles or author or any other combination of fields that identify the document.

You specify the name of structured column or columns in the SELECT statement.

Ordering the Hitlist

When you issue either a text query or theme query, Oracle returns the hitlist of documents that satisfy the query with a relevance score for each document returned. You can use these scores to order the hitlist to show the most relevant documents first.

The score for each document is between one and one hundred. The higher the score, the more relevant the document.

Oracle calculates scores when you use the CONTAINS and CATSEARCH operator. You obtain scores using the SCORE operator.

See Also: [Chapter 3, "Querying"](#)

Presenting Document Hit Count

You present the number of hits the query returned alongside the hitlist, using SELECT COUNT(*). For example:

```
SELECT COUNT(*) FROM docs WHERE CONTAINS(text, 'oracle', 1) > 0;
```

To count hits in PL/SQL, you can also use the CTX_QUERY.COUNT_HITS procedure.

Document Presentation and Highlighting

Typically, a query application allows the user to view the documents returned by a query. The user selects a document from the hitlist and then the application presents the document in some form.

With Oracle Text, you can display a document in different ways. For example, you can present documents with query terms highlighted. Highlighted query terms can be either the words of a word query or the themes of an ABOUT query in English.

You can also obtain gist (document summary) and theme information from documents with the CTX_DOC PL/SQL package.

[Table 1-1](#) describes the different output you can obtain and which procedure to use to obtain each type:

Table 1-1

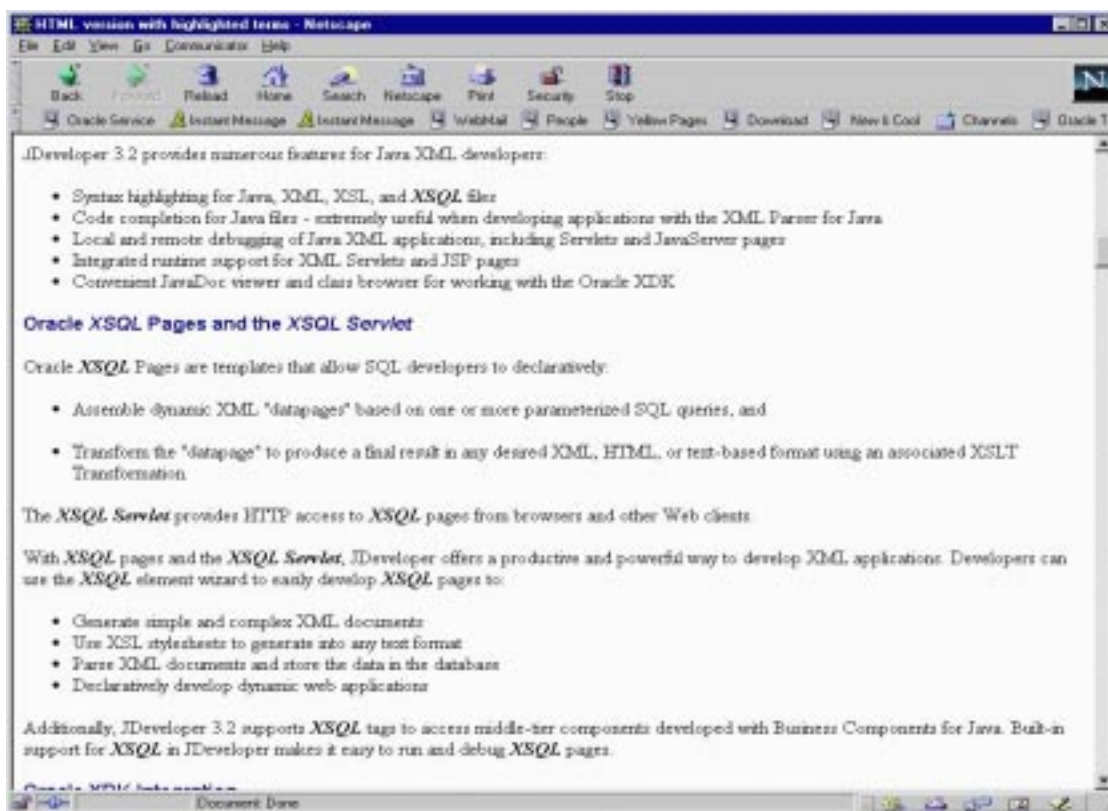
Output	Procedure
Plain text version, no highlights	CTX_DOC.FILTER
HTML version of document, no highlights	CTX_DOC.FILTER
Highlighted document, plain text version	CTX_DOC.MARKUP
Highlighted document, HTML version	CTX_DOC.MARKUP
Highlight offset information for plain text version	CTX_DOC.HIGHLIGHT
Highlight offset information for HTML version	CTX_DOC.HIGHLIGHT
Theme summaries and gist of document.	CTX_DOC.GIST
List of themes in document.	CTX_DOC.THEMES

See Also: [Chapter 4, "Document Presentation"](#)

Highlighting Example

Figure 1–5 is a screen shot of a query application presenting a document with the query terms *XSQL* and *Servlet* highlighted.

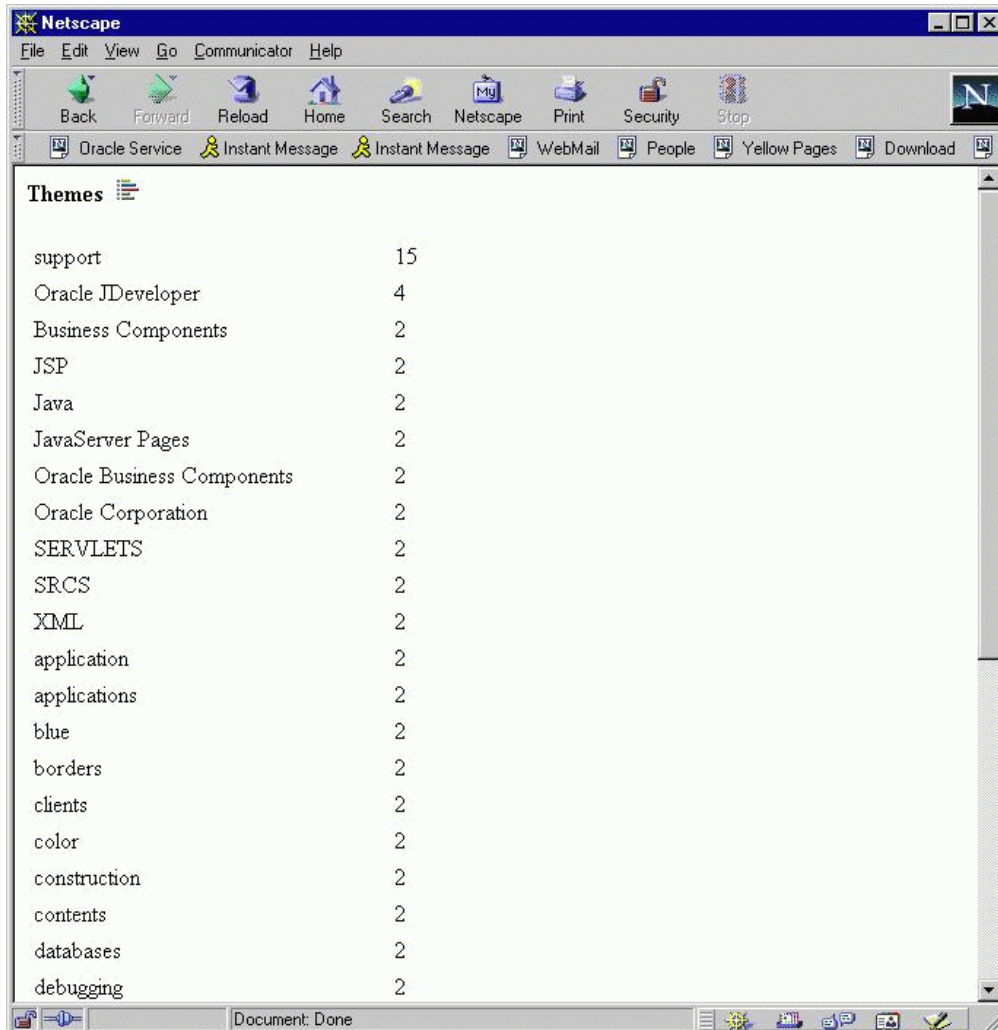
Figure 1–5 Query Application Presenting Highlighted Document



Document List of Themes Example

Figure 1-6 is a screen shot of a query application presenting a list of themes for a document.

Figure 1-6 Query Application Displaying Document Themes



Gist Example

Figure 1-7 is a screen shot of a query application presenting a gist for a document.

Figure 1-7 Query Application Presenting Document Gist



The chapter is an introduction to Oracle Text indexing. The following topics are covered:

- [About Oracle Text Indexes](#)
- [Considerations For Indexing](#)
- [Index Creation](#)
- [Index Maintenance](#)
- [Managing DML Operations for a CONTEXT Index](#)

About Oracle Text Indexes

An Oracle Text index is an Oracle domain index. To build your query application, you can create an Oracle Text index of type `CONTEXT` and query it with the `CONTAINS` operator.

For better performance for mixed queries, you can create a `CTXCAT` index. Use this index type when your application relies heavily on mixed queries to search small documents or descriptive text fragments based on related criteria such as dates or prices. You query this index with the `CATSEARCH` operator.

To build a document classification application, you create an Oracle Text index of type `CTXRULE`. With such an index, you can classify plain text, HTML, or XML documents using the `MATCHES` operator.

You create an index from a populated text table. In a query application, the table must contain the text or pointers to where the text is stored. Text is usually a collection of documents, but can also be small text fragments. If you are building a document classification application, you store your defining query set in the text table.

You create a text index as a type of extensible index to Oracle using standard SQL. This means that an Oracle Text index operates like an Oracle index. It has a name by which it is referenced and can be manipulated with standard SQL statements.

The benefits of a creating an Oracle Text index include fast response time for text queries with the `CONTAINS`, `CATSEARCH`, and `MATCHES` Oracle Text operators. These operators query the `CONTEXT`, `CTXCAT`, and `CTXRULE` index types respectively.

See Also: For more information about creating a Text index, see ["Index Creation"](#) in this chapter.

Structure of the Oracle Text `CONTEXT` Index

Oracle Text indexes text by converting all words into tokens. The general structure of an Oracle Text `CONTEXT` index is an inverted index where each token contains the list of documents (rows) that contain that token.

For example, after a single initial indexing operation, the word `DOG` might have an entry as follows:

```
DOG DOC1 DOC3 DOC5
```

This means that the word DOG is contained in the rows that store documents one, three and five.

For more information, see [optimizing the index](#) in this chapter.

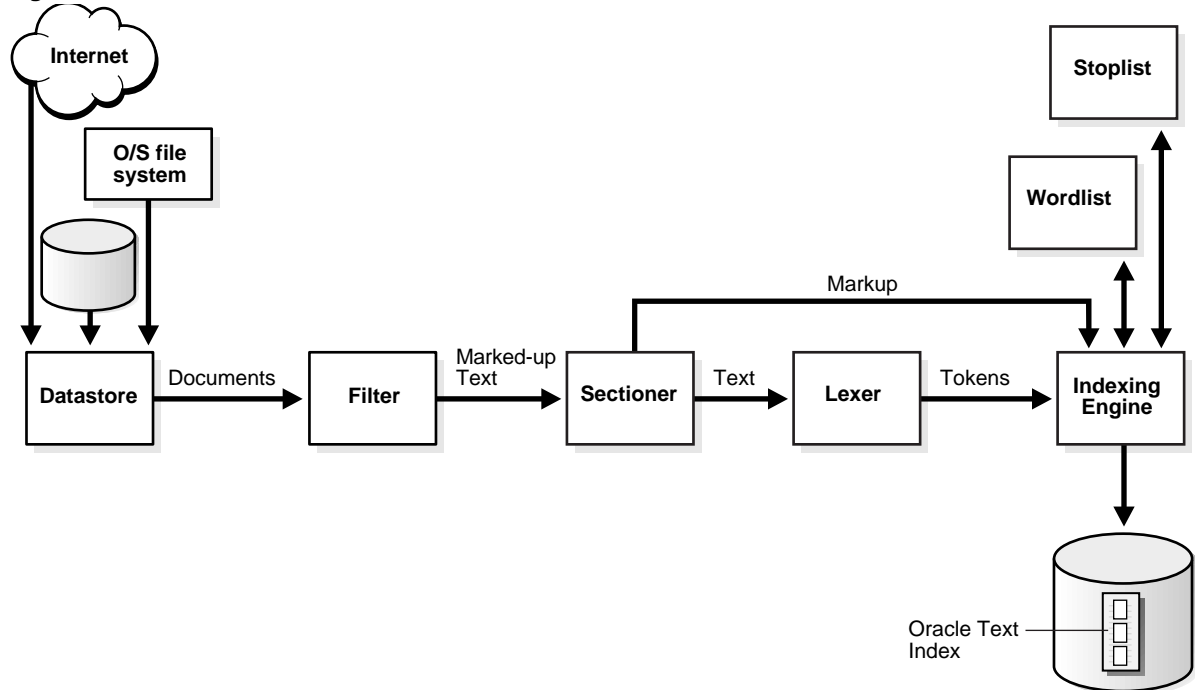
Merged Word and Theme Index

By default in English and French, Oracle Text indexes theme information with word information. You can query theme information with the ABOUT operator. You can optionally enable and disable theme indexing.

See Also: To learn more about indexing theme information, see ["Creating Preferences"](#) in this chapter.

The Oracle Text Indexing Process

Figure 2-1



You initiate the indexing process with the `CREATE INDEX` statement. The goal is to create an Oracle Text index of tokens according to the parameters and preferences you specify.

Figure 2-1 shows the indexing process. This process is a data stream that is acted upon by the different indexing objects. Each object corresponds to an indexing preference type or section group you can specify in the parameter string of `CREATE INDEX` or `ALTER INDEX`. The sections that follow describe these objects.

Datastore Object

The stream starts with the datastore reading in the documents as they are stored in the system according to your datastore preference. For example, if you have defined your datastore as `FILE_DATASTORE`, the stream starts by reading the files from the operating system. You can also store your documents on the internet or in the Oracle database.

Filter Object

The stream then passes through the filter. What happens here is determined by your FILTER preference. The stream can be acted upon in one of the following ways:

- No filtering takes place. This happens when you specify the NULL_FILTER preference type. Documents that are plain text, HTML, or XML need no filtering.
- Formatted documents (binary) are filtered to marked-up text. This happens when you specify the INSO_FILTER preference type.
- Text is converted from a non-database character set to the database character set. This happens when you specify CHARSET_FILTER preference type.

Sectioner Object

After being filtered, the marked-up text passes through the sectioner that separates the stream into text and section information. Section information includes where sections begin and end in the text stream. The type of sections extracted is determined by your section group type.

The section information is passed directly to the indexing engine which uses it later. The text is passed to the lexer.

Lexer Object

The lexer breaks the text into tokens according to your language. These tokens are usually words. To extract tokens, the lexer uses the parameters as defined in your lexer preference. These parameters include the definitions for the characters that separate tokens such as whitespace, and whether to convert the text to all uppercase or to leave it in mixed case.

When theme indexing is enabled, the lexer analyses your text to create theme tokens for indexing.

Indexing Engine

The indexing engine creates the inverted index that maps tokens to the documents that contain them. In this phase, Oracle uses the stoplist you specify to exclude stopwords or stopthemes from the index. Oracle also uses the parameters defined in your WORDLIST preference, which tell the system how to create a prefix index or substring index, if enabled.

Partitioned Tables and Indexes

You can create a partitioned CONTEXT index on a partitioned text table. The table must be partitioned by range. Hash, composite and list partitions are not supported.

You might create a partitioned text table to partition your data by date. For example, if your application maintains a large library of dated news articles, you can partition your information by month or year. Partitioning simplifies the manageability of large databases since querying, DML, and backup and recovery can act on single partitions.

See Also: *Oracle9i Database Concepts* for more information about partitioning.

Querying Partitioned Tables

To query a partitioned table, you use CONTAINS in the SELECT statement no differently as you query a regular table. You can query the entire table or a single partition. However, if you are using the ORDER BY SCORE clause, Oracle recommends that you query single partitions unless you include a range predicate that limits the query to a single partition.

Parallel Indexing

Oracle Text supports parallel indexing with CREATE INDEX on a partitioned text table.

The parallel indexing operation creates multiple threads where each thread works on a partition. Since indexing is an I/O intensive operation, parallel indexing is most effective in decreasing your indexing time when you have distributed disk access and multiple CPUs.

Since parallel indexing decreases the *initial* indexing time, it is useful for

- data staging, when your product includes an Oracle Text index
- rapid initial startup of applications based on large data collections
- application testing, when you need to test different index parameters and schemas while developing your application

Note: Parallel indexing with a partitioned text table can only affect the performance of an initial index with CREATE INDEX. It does not affect DML performance with ALTER INDEX, and has minimal impact on query performance.

Limitations for Indexing

Columns with Multiple Indexes

A column can have no more than a single domain index attached to it, which is in keeping with Oracle standards. However, a single Text index can contain theme information in addition to word information.

Indexing Views

Oracle SQL standards does not support creating indexes on views. Therefore, if you need to index documents whose contents are in different tables, you can create a data storage preference using the USER_DATASTORE object. With this object, you can define a procedure that synthesizes documents from different tables at index time.

See Also: *Oracle Text Reference* to learn more about USER_DATASTORE.

Considerations For Indexing

You use the `CREATE INDEX` statement to create an Oracle Text index. When you create an index and specify no parameter string, an index is created with default parameters.

You can also override the defaults and customize your index to suit your query application. The parameters and preference types you use to customize your index with `CREATE INDEX` fall into the following general categories.

Type of Index

With Oracle Text, you can create one of three index types with CREATE INDEX. The following table describes each type, its purpose, and what features it supports:

Index Type	Description	Supported Preferences and Parameters	Query Operator	Notes
CONTEXT	<p>Use this index to build a text retrieval application when your text consists of large coherent documents. You can index documents of different formats such as MS Word, HTML or plain text.</p> <p>With a context index, you can customize your index in a variety of ways.</p>	<p>All CREATE INDEX preferences and parameters supported except for INDEX SET.</p> <p>These supported parameters include the index partition clause, and the format, charset, and language columns.</p>	CONTAINS	<p>Supports all documents services and query services.</p> <p>Supports indexing of partitioned text tables.</p>
CTXCAT	<p>Use this index type for better mixed query performance. Typically, with this index type, you index small documents or text fragments. Other columns in the base table, such as item names, prices and descriptions can be included in the index to improve mixed query performance.</p>	<p>INDEX SET</p> <p>LEXER (theme indexing not supported)</p> <p>STOPLIST</p> <p>STORAGE</p> <p>WORDLIST (only prefix_index attribute supported for Japanese data)</p> <p>Format, charset, and language columns not supported.</p> <p>Table and index partitioning not supported.</p>	<p>CATSEARCH</p> <p>This operator has its own query language that supports logical operations, phrase queries, and wildcarding.</p> <p>The query language does not support ABOUT, fuzzy, and stem operators.</p>	<p>The size of a CTXCAT index is related to the total amount of text to be indexed, number of indexes in the index set, and number of columns indexed. Carefully consider your queries and your resources before adding indexes to the index set.</p> <p>The CTXCAT index does not support table and index partitioning, documents services (highlighting, markup, themes, and gists) or query services (explain, query feedback, and browse words.)</p>

Index Type	Description	Supported Preferences and Parameters	Query Operator	Notes
CTXRULE	Use CTRRULE index to build a document classification or routing application. The CTRRULE index is an index created on a table of queries, where the queries define the classification or routing criteria.	<p>Only the BASIC_LEXER type supported for indexing your query set.</p> <p>Queries in your query set can include ABOUT, STEM, AND, NEAR, NOT, and OR operators.</p> <p>The following operators are not supported: ACCUM, EQUIV, WITHIN, WILDCARD, FUZZY, SOUNDEX, MINUS, WEIGHT, THRESHOLD.</p> <p>The CREATE INDEX storage clause supported for creating the index on the queries.</p> <p>Section group supported for when you use the MATCHES operator to classify documents.</p> <p>Wordlist supported for stemming operations on your query set.</p> <p>Filter, memory, datastore, and populate parameters are not applicable to index type CTRRULE.</p>	MATCHES	Single documents (plain text, HTML, or XML) can be classified using the MATCHES operator, which turns a document into a set of queries and finds the matching rows in the CTRRULE index.

See Also: [Index Creation](#) in this chapter.

Location of Text

Your document text can reside in one of three places, the text table, the file system, or the world-wide web. When you index with CREATE INDEX, you specify the location using the datastore preference. Use the appropriate datastore according to your application.

The following table describes all the different ways you can store your text with the datastore preference type.

Datastore Type	Use When
DIRECT_DATASTORE	Data is stored internally in a text column. Each row is indexed as a single document. Your text column can be VARCHAR2, CLOB, BLOB, CHAR, or BFILE. XMLType columns are supported for the context index type.
MULTI_COLUMN_DATASTORE	Data is stored in a text table in more than one column. Columns are concatenated to create a virtual document, one document per row.
DETAIL_DATASTORE	Data is stored internally in a text column. Document consists of one or more rows stored in a text column in a detail table, with header information stored in a master table.
FILE_DATASTORE	Data is stored externally in operating system files. Filenames are stored in the text column, one per row.
NESTED_DATASTORE	Data is stored in a nested table.
URL_DATASTORE	Data is stored externally in files located on an intranet or the Internet. Uniform Resource Locators (URLs) are stored in the text column.
USER_DATASTORE	Documents are synthesized at index time by a user-defined stored procedure.

Indexing time and document retrieval time will be increased for indexing URLs since the system must retrieve the document from the network.

See Also: [Datastore Examples](#) in this chapter.

Document Formats and Filtering

Formatted documents such as Microsoft Word and PDF must be filtered to text to be indexed. The type of filtering the system uses is determined by the FILTER preference type. By default the system uses the INSO_FILTER filter type which automatically detects the format of your documents and filters them to text.

Oracle can index most formats. Oracle can also index columns that contain documents with mixed formats.

No Filtering for HTML

If you are indexing HTML or plain text files, do not use the `INSO_FILTER` type. For best results, use the `NULL_FILTER` preference type.

See Also: [NULL_FILTER Example: Indexing HTML Documents](#) in this chapter.

Filtering Mixed Formatted Columns

If you have a mixed format column such as one that contains Microsoft Word, plain text, and HTML documents, you can bypass filtering for plain text or HTML by including a format column in your text table. In the format column, you tag each row `TEXT` or `BINARY`. Rows that are tagged `TEXT` are not filtered.

For example, you can tag the HTML and plain text rows as `TEXT` and the Microsoft Word rows as `BINARY`. You specify the format column in the `CREATE INDEX` parameter clause.

Custom Filtering

You can create your own custom filter to filter documents for indexing. You can create either an external filter that is executed from the file system or an internal filter as a PL/SQL or Java stored procedure.

For external custom filtering, use the `USER_FILTER` filter preference type.

For internal filtering, use the `PROCEDURE_FILTER` filter type.

See Also: [PROCEDURE_FILTER Example](#) in this chapter.

Bypassing Rows for Indexing

You can bypass rows in your text table that are not to be indexed, such as rows that contain image data. To do so, create a format column in your table and set it to `IGNORE`. You name the format column in the parameter clause of `CREATE INDEX`.

Document Character Set

The indexing engine expects filtered text to be in the database character set. When you use the `INSO_FILTER` filter type, formatted documents are converted to text in the database character set.

If your source is text and your document character set is not the database character set, you can use the `INSO_FILTER` or `CHARSET_FILTER` filter type to convert your text for indexing.

Mixed Character Set Columns

If your document set contains documents with different character sets, such as `JA16EUC` and `JA16SJIS`, you can index the documents provided you create a `charset` column. You populate this column with the name of the document character set on a per-row basis. You name the column in the parameter clause of the `CREATE INDEX` statement.

Document Language

Oracle can index most languages. By default, Oracle assumes the language of text to index is the language you specify in your database setup.

You use the `BASIC_LEXER` preference type to index whitespace-delimited languages such as English, French, German, and Spanish. For some of these languages you can enable alternate spelling, composite word indexing, and base letter conversion.

You can also index Japanese, Chinese, and Korean.

See Also: *Oracle Text Reference* to learn more about indexing these languages.

Indexing Multi-language Columns

Oracle can index text columns that contain documents of different languages, such as a column that contains documents written in English, German, and Japanese. To index a multi-language column, you need a language column in your text table. Use the `MULTI_LEXER` preference type.

You can also incorporate a multi-language stoplist when you index multi-language columns.

See Also: [MULTI_LEXER Example: Indexing a Multi-Language Table](#) in this chapter.

Indexing Special Characters

When you use the `BASIC_LEXER` preference type, you can specify how non-alphanumeric characters such as hyphens and periods are indexed with respect

to the tokens that contain them. For example, you can specify that Oracle include or exclude hyphen character (-) when indexing a word such as *web-site*.

These characters fall into BASIC_LEXER categories according to the behavior you require during indexing. The way the you set the lexer to behave for indexing is the way it behaves for query parsing.

Some of the special characters you can set are as follows:

Printjoins Character

Define a non-alphanumeric character as printjoin when you want this character to be included in the token during indexing.

For example, if you want your index to include hyphens and underscore characters, define them as printjoins. This means that words such as *web-site* are indexed as *web-site*. A query on *website* does not find *web-site*.

See Also: [BASIC_LEXER Example: Setting Printjoins Characters](#) in this chapter.

Skipjoins Character

Define a non-alphanumeric character as a skipjoin when you do not want this character to be indexed with the token that contains it.

For example, with the hyphen (-) character defined as a skipjoin, the word *web-site* is indexed as *website*. A query on *web-site* finds documents containing *website* and *web-site*.

Other Characters

Other characters can be specified to control other tokenization behavior such as token separation (startjoins, endjoins, whitespace), punctuation identification (punctuations), number tokenization (numjoins), and word continuation after line-breaks (continuation). These categories of characters have defaults, which you can modify.

See Also: *Oracle Text Reference* to learn more about the BASIC_LEXER.

Case-Sensitive Indexing and Querying

By default, all text tokens are converted to uppercase and then indexed. This results in case-insensitive queries. For example, separate queries on each of the three words *cat*, *CAT*, and *Cat* all return the same documents.

You can change the default and have the index record tokens as they appear in the text. When you create a case-sensitive index, you must specify your queries with exact case to match documents. For example, if a document contains *Cat*, you must specify your query as *Cat* to match this document. Specifying *cat* or *CAT* does not return the document.

To enable or disable case-sensitive indexing, use the `mixed_case` attribute of the `BASIC_LEXER` preference.

See Also: *Oracle Text Reference* to learn more about the `BASIC_LEXER`.

Language Specific Features

You can enable the following language specific features at index time:

Indexing Themes

For English and French, you can index document theme information. A document theme is a main document concept. Themes can be queried with the `ABOUT` operator.

You can index theme information in other languages provided you have loaded and compiled a knowledge base for the language.

By default themes are indexed in English and French. You can enable and disable theme indexing with the `index_themes` attribute of the `BASIC_LEXER` preference type.

See Also: *Oracle Text Reference* to learn more about the `BASIC_LEXER`.

[ABOUT Queries and Themes](#) in [Chapter 3, "Querying"](#).

Base-Letter Conversion for Characters with Diacritical Marks

Some languages contain characters with diacritical marks such as tildes, umlauts, and accents. When your indexing operation converts words containing diacritical

marks to their base letter form, queries need not contain diacritical marks to score matches. For example in Spanish with a base-letter index, a query of *energía* matches *energía* and *energía* in the index.

However, with base-letter indexing disabled, a query of *energía* matches only *energía*.

You can enable and disable base-letter indexing for your language with the `base_` letter attribute of the `BASIC_LEXER` preference type.

See Also: *Oracle Text Reference* to learn more about the `BASIC_LEXER`.

Alternate Spelling

Languages such as German, Danish, and Swedish contain words that have more than one accepted spelling. For instance, in German, the *ä* character can be substituted for the *ae* character. The *ae* character is known as the base letter form.

By default, Oracle indexes words in their base-letter form for these languages. Query terms are also converted to their base-letter form. The result is that these words can be queried with either spelling.

You can enable and disable alternate spelling for your language using the `alternate_` spelling attribute in the `BASIC_LEXER` preference type.

See Also: *Oracle Text Reference* to learn more about the `BASIC_LEXER`.

Composite Words

German and Dutch text contain composite words. By default, Oracle creates composite indexes for these languages. The result is that a query on a term returns words that contain the term as a sub-composite.

For example, in German, a query on the term *Bahnhof* (train station) returns documents that contain *Bahnhof* or any word containing *Bahnhof* as a sub-composite, such as *Hauptbahnhof*, *Nordbahnhof*, or *Ostbahnhof*.

You can enable and disable the creation of composite indexes with the `composite` attribute of the `BASIC_LEXER` preference.

See Also: *Oracle Text Reference* to learn more about the `BASIC_LEXER`.

Korean, Japanese, and Chinese Indexing

You index these languages with specific lexers:

Language	Lexer
Korean	KOREAN_MORPH_LEXER
Japanese	JAPANESE_LEXER
Chinese	CHINESE_VGRAM_LEXER

The KOREAN_MORPH_LEXER has its own set of attributes to control indexing. Features include composite word indexing.

See Also: *Oracle Text Reference* to learn more about these lexers.

Fuzzy Matching and Stemming

Fuzzy matching enables you to match similarly spelled words in queries. Stemming enables you to match words with the same linguistic root.

Fuzzy matching and stemming are automatically enabled in your index if Oracle Text supports this feature for your language.

Fuzzy matching is enabled with default parameters for its similarity score lower limit and for its maximum number of expanded terms. At index time you can change these default parameters.

See Also: *Oracle Text Reference* for more information about the BASIC_WORDLIST preference type.

Better Wildcard Query Performance

Wildcard queries enable you to issue left-truncated, right-truncated and doubly truncated queries, such as *%ing*, *cos%*, or *%benz%*. With normal indexing, these queries can sometimes expand into large word lists, degrading your query performance.

Wildcard queries have better response time when token prefixes and substrings are recorded in the index.

By default, token prefixes and substrings are not recorded in the Oracle Text index. If your query application makes heavy use of wildcard queries, consider indexing token prefixes and substrings. To do so, use the wordlist preference type. The trade-off is a bigger index for improved wildcard searching.

See Also: [BASIC_WORDLIST Example: Enabling Substring and Prefix Indexing](#) in this chapter.

Document Section Searching

For documents that have internal structure such as HTML and XML, you can define and index document sections. Indexing document sections enables you to narrow the scope of your queries to within pre-defined sections. For example, you can specify a query to find all documents that contain the term *dog* within a section you define as *Headings*.

Sections must be defined prior to indexing and specified with the section group preference.

Oracle Text provides section groups with system-defined section definitions for HTML and XML. You can also specify that the system automatically create sections from XML documents during indexing.

See Also: [Chapter 6, "Document Section Searching"](#)

Stopwords and Stopthemes

A stopword is a word that is not to be indexed. Usually stopwords are low information words in a given language such as *this* and *that* in English.

By default, Oracle provides a list of stopwords called a stoplist for indexing a given language. You can modify this list or create your own with the CTX_DDL package. You specify the stoplist in the parameter string of CREATE INDEX.

A stoptheme is a word that is prevented from being theme-indexed or prevented from contributing to a theme. You can add stopthemes with the CTX_DDL package.

You can search document themes with the ABOUT operator. You can retrieve document themes programatically with the CTX_DOC PL/SQL package.

Multi-Language Stoplists

You can also create multi-language stoplists to hold language-specific stopwords. A multi-language stoplist is useful when you use the MULTI_LEXER to index a table that contains documents in different languages, such as English, German, and Japanese.

At indexing time, the language column of each document is examined, and only the stopwords for that language are eliminated. At query time, the session language

setting determines the active stopwords, like it determines the active lexer when using the multi-lexer.

Index Creation

You can create three types of indexes with Oracle Text: CONTEXT, CTXCAT, and CTXRULE.

Procedure for Creating a CONTEXT Index

By default, the system expects your documents to be stored in a text column. Once this requirement is satisfied, you can create a text index using the CREATE INDEX SQL command as an extensible index of type context, without explicitly specifying any preferences. The system automatically detects your language, the datatype of the text column, format of documents, and sets indexing preferences accordingly.

See Also: For more information about the out-of-box defaults, see [Default CONTEXT Index Example](#) in this chapter.

To create an Oracle Text index, do the following:

1. Optionally, determine your custom indexing preferences, section groups, or stoplists if not using defaults. The following table describes these indexing classes:

Class	Description
Datastore	How are your documents stored?
Filter	How can the documents be converted to plaintext?
Lexer	What language is being indexed?
Wordlist	How should stem and fuzzy queries be expanded?
Storage	How should the index data be stored?
Stop List	What words or themes are not to be indexed?
Section Group	How are documents sections defined?

See Also: [Considerations For Indexing](#) in this chapter and *Oracle Text Reference*.

2. Optionally, create your own custom preferences, section groups, or stoplists. See ["Creating Preferences"](#) in this chapter.

3. Create the Text index with the SQL command CREATE INDEX, naming your index and optionally specifying preferences. See "[Creating an Index](#)" in this chapter.

Creating Preferences

You can optionally create your own custom index preferences to override the defaults. Use the preferences to specify index information such as where your files are stored and how to filter your documents. You create the preferences then set the attributes.

Datastore Examples

The following sections give examples for setting direct, multi-column, URL, and file datastores.

See Also: *Oracle Text Reference* for more information about data storage.

Specifying DIRECT_DATASTORE The following example creates a table with a CLOB column to store text data. It then populates two rows with text data and indexes the table using the system-defined preference CTXSYS.DEFAULT_DATASTORE.

```
create table mytable(id number primary key, docs clob);

insert into mytable values(111555,'this text will be indexed');
insert into mytable values(111556,'this is a direct_datastore example');
commit;

create index myindex on mytable(docs)
  indextype is ctxsys.context
  parameters ('DATASTORE CTXSYS.DEFAULT_DATASTORE');
```

Specifying MULTI_COLUMN_DATASTORE The following example creates a multi-column datastore preference called `my_multi` on the three text columns to be concatenated and indexed:

```
begin
ctx_ddl.create_preference('my_multi', 'MULTI_COLUMN_DATASTORE');
ctx_ddl.set_attribute('my_multi', 'columns', 'column1, column2, column3');
end;
```

Specifying URL Data Storage This example creates a URL_DATASTORE preference called my_url to which the http_proxy, no_proxy, and timeout attributes are set. The defaults are used for the attributes that are not set.

```
begin
  ctx_ddl.create_preference('my_url', 'URL_DATASTORE');
  ctx_ddl.set_attribute('my_url', 'HTTP_PROXY', 'www-proxy.us.oracle.com');
  ctx_ddl.set_attribute('my_url', 'NO_PROXY', 'us.oracle.com');
  ctx_ddl.set_attribute('my_url', 'Timeout', '300');
end;
```

Specifying File Data Storage The following example creates a data storage preference using the FILE_DATASTORE. This tells the system that the files to be indexed are stored in the operating system. The example uses CTX_DDL.SET_ATTRIBUTE to set the PATH attribute of to the directory /docs.

```
begin
  ctx_ddl.create_preference('mypref', 'FILE_DATASTORE');
  ctx_ddl.set_attribute('mypref', 'PATH', '/docs');
end;
```

NULL_FILTER Example: Indexing HTML Documents

If your document set is entirely HTML, Oracle recommends that you use the NULL_FILTER in your filter preference, which does no filtering.

For example, to index an HTML document set, you can specify the system-defined preferences for NULL_FILTER and HTML_SECTION_GROUP as follows:

```
create index myindex on docs(htmlfile) indextype is ctxsys.context
  parameters('filter ctxsys.null_filter
  section group ctxsys.html_section_group');
```

PROCEDURE_FILTER Example

Consider a filter procedure CTXSYS.NORMALIZE that you define with the following signature:

```
PROCEDURE NORMALIZE(id IN ROWID, charset IN VARCHAR2, input IN CLOB,
  output IN OUT NOCOPY VARCHAR2);
```

To use this procedure as your filter, you set up your filter preference as follows:

```
begin
  ctx_ddl.create_preference('myfilt', 'procedure_filter');
```

```

    ctx_ddl.set_attribute('myfilt', 'procedure', 'normalize');
    ctx_ddl.set_attribute('myfilt', 'input_type', 'clob');
    ctx_ddl.set_attribute('myfilt', 'output_type', 'varchar2');
    ctx_ddl.set_attribute('myfilt', 'rowid_parameter', 'TRUE');
    ctx_ddl.set_attribute('myfilt', 'charset_parameter', 'TRUE');
end;

```

BASIC_LEXER Example: Setting Printjoins Characters

Printjoin characters are non-alphanumeric characters that are to be included in index tokens, so that words such as *web-site* are indexed as *web-site*.

The following example sets printjoin characters to be the hyphen and underscore with the BASIC_LEXER:

```

begin
    ctx_ddl.create_preference('mylex', 'BASIC_LEXER');
    ctx_ddl.set_attribute('mylex', 'printjoins', '-_');
end;

```

To create the index with printjoins characters set as above, issue the following statement:

```

create index myindex on mytable ( docs )
    indextype is ctxsys.context
    parameters ( 'LEXER mylex' );

```

MULTI_LEXER Example: Indexing a Multi-Language Table

You use the MULTI_LEXER preference type to index a column containing documents in different languages. For example, you can use this preference type when your text column stores documents in English, German, and French.

The first step is to create the multi-language table with a primary key, a text column, and a language column as follows:

```

create table globaldoc (
    doc_id number primary key,
    lang varchar2(3),
    text clob
);

```

Assume that the table holds mostly English documents, with some German and Japanese documents. To handle the three languages, you must create three sub-lexers, one for English, one for German, and one for Japanese:

```

ctx_ddl.create_preference('english_lexer', 'basic_lexer');

```

```
ctx_ddl.set_attribute('english_lexer','index_themes','yes');
ctx_ddl.set_attribute('english_lexer','theme_language','english');

ctx_ddl.create_preference('german_lexer','basic_lexer');
ctx_ddl.set_attribute('german_lexer','composite','german');
ctx_ddl.set_attribute('german_lexer','mixed_case','yes');
ctx_ddl.set_attribute('german_lexer','alternate_spelling','german');

ctx_ddl.create_preference('japanese_lexer','japanese_vgram_lexer');
```

Create the multi-lexer preference:

```
ctx_ddl.create_preference('global_lexer','multi_lexer');
```

Since the stored documents are mostly English, make the English lexer the default using `CTX_DDL.ADD_SUB_LEXER`:

```
ctx_ddl.add_sub_lexer('global_lexer','default','english_lexer');
```

Now add the German and Japanese lexers in their respective languages with `CTX_DDL.ADD_SUB_LEXER` procedure. Also assume that the language column is expressed in the standard ISO 639-2 language codes, so add those as alternate values.

```
ctx_ddl.add_sub_lexer('global_lexer','german','german_lexer','ger');
ctx_ddl.add_sub_lexer('global_lexer','japanese','japanese_lexer','jpn');
```

Now create the index `globalx`, specifying the multi-lexer preference and the language column in the parameter clause as follows:

```
create index globalx on globaldoc(text) indextype is ctxsys.context
parameters ('lexer global_lexer language column lang');
```

BASIC_WORDLIST Example: Enabling Substring and Prefix Indexing

The following example sets the wordlist preference for prefix and substring indexing. Having a prefix and sub-string component to your index improves performance for wildcard queries.

For prefix indexing, the example specifies that Oracle create token prefixes between three and four characters long:

```
begin
  ctx_ddl.create_preference('mywordlist','BASIC_WORDLIST');
  ctx_ddl.set_attribute('mywordlist','INDEX_PREFIX','YES');
  ctx_ddl.set_attribute('mywordlist','PREFIX_MIN_LENGTH',3);
  ctx_ddl.set_attribute('mywordlist','PREFIX_MAX_LENGTH',4);
end;
```

```
        ctx_ddl.set_attribute('mywordlist','SUBSTRING_INDEX','YES');
    end
```

Creating Section Groups for Section Searching

When documents have internal structure such as in HTML and XML, you can define document sections using embedded tags before you index. This enables you to query within the sections using the WITHIN operator. You define sections as part of a section group.

Example: Creating HTML Sections

The following code defines a section group called `htmgroup` of type `HTML_SECTION_GROUP`. It then creates a zone section in `htmgroup` called `heading` identified by the `<H1>` tag:

```
begin
ctx_ddl.create_section_group('htmgroup', 'HTML_SECTION_GROUP');
ctx_ddl.add_zone_section('htmgroup', 'heading', 'H1');
end;
```

See Also: [Chapter 6, "Document Section Searching"](#)

Using Stopwords and Stoplists

A stopword is a word that is not to be indexed. A stopword is usually a low information word such as *this* or *that* in English.

The system supplies a list of stopwords called a stoplist for every language. By default during indexing, the system uses the Oracle Text default stoplist for your language.

You can edit the default stoplist `CTXSYS.DEFAULT_STOPLIST` or create your own with the following PL/SQL procedures:

- `CTX_DDL.CREATE_STOPLIST`
- `CTX_DDL.ADD_STOPWORD`
- `CTX_DDL.REMOVE_STOPWORD`

You specify your custom stoplists in the parameter clause of `CREATE INDEX`.

You can also dynamically add stopwords after indexing with the `ALTER INDEX` statement.

Multi-Language Stoplists

You can create multi-language stoplists to hold language-specific stopwords. A multi-language stoplist is useful when you use the `MULTI_LEXER` to index a table that contains documents in different languages, such as English, German, and Japanese.

To create a multi-language stoplist, use the `CTX_DLL.CREATE_STOPLIST` procedure and specify a stoplist type of `MULTI_STOPLIST`. You add language specific stopwords with `CTX_DDL.ADD_STOPWORD`.

Stopthemes and Stopclasses

In addition to defining your own stopwords, you can define stopthemes, which are themes that are not to be indexed. This feature is available for English only.

You can also specify that numbers are not to be indexed. A class of alphanumeric characters such a numbers that is not to be indexed is a *stopclass*.

You record your own stopwords, stopthemes, stopclasses by creating a single stoplist, to which you add the stopwords, stopthemes, and stopclasses. You specify the stoplist in the paramstring for `CREATE INDEX`.

PL/SQL Procedures for Managing Stoplists

You use the following procedures to manage stoplists, stopwords, stopthemes, and stopclasses:

- `CTX_DDL.CREATE_STOPLIST`
- `CTX_DDL.ADD_STOPWORD`
- `CTX_DDL.ADD_STOPTHEME`
- `CTX_DDL.ADD_STOPCLASS`
- `CTX_DDL.REMOVE_STOPWORD`
- `CTX_DDL.REMOVE_STOPTHEME`
- `CTX_DDL.REMOVE_STOPCLASS`
- `CTX_DDL.DROP_STOPLIST`

See Also: *Oracle Text Reference*, to learn more about using these commands.

Creating an Index

You create an Oracle Text index as an extensible index using the CREATE INDEX SQL command.

You can create three types of indexes:

- CONTEXT
- CTXCAT
- CTXRULE

Creating a CONTEXT Index

The context index type is well-suited for indexing large coherent documents such as MS Word, HTML or plain text. With a context index, you can also customize your index in a variety of ways.

The documents must be loaded in a text table.

Default CONTEXT Index Example

The following command creates a default context index called `myindex` on the `text` column in the `docs` table:

```
CREATE INDEX myindex ON docs(text) INDEXTYPE IS CTXSYS.CONTEXT;
```

When you use CREATE INDEX without explicitly specifying parameters, the system does the following for all languages by default:

- Assumes that the text to be indexed is stored directly in a text column. The text column can be of type CLOB, BLOB, BFILE, VARCHAR2, or CHAR.
- Detects the column type and uses filtering for binary column types. Most document formats are supported for filtering. If your column is plain text, the system does not use filtering.

Note: For document filtering to work correctly in your system, you must ensure that your environment is set up correctly to support the Inso filter.

To learn more about configuring your environment to use the Inso filter, see the *Oracle Text Reference*.

- Assumes the language of text to index is the language you specify in your database setup.

- Uses the default stoplist for the language you specify in your database setup. Stoplists identify the words that the system ignores during indexing.
- Enables fuzzy and stemming queries for your language, if this feature is available for your language.

You can always change the default indexing behavior by creating your own preferences and specifying these custom preferences in the parameter string of `CREATE INDEX`.

Custom `CONTEXT` Index Example: Indexing HTML Documents

To index an HTML document set located by URLs, you can specify the system-defined preference for the `NULL_FILTER` in the `CREATE INDEX` statement.

You can also specify your section group `htmgroup` that uses `HTML_SECTION_GROUP` and datastore `my_url` that uses `URL_DATASTORE` as follows:

```
begin
  ctx_ddl.create_preference('my_url', 'URL_DATASTORE');
  ctx_ddl.set_attribute('my_url', 'HTTP_PROXY', 'www-proxy.us.oracle.com');
  ctx_ddl.set_attribute('my_url', 'NO_PROXY', 'us.oracle.com');
  ctx_ddl.set_attribute('my_url', 'Timeout', '300');
end;
```

```
begin
  ctx_ddl.create_section_group('htmgroup', 'HTML_SECTION_GROUP');
  ctx_ddl.add_zone_section('htmgroup', 'heading', 'H1');
end;
```

You can then index your documents as follows:

```
create index myindex on docs(htmlfile) indextype is ctxsys.context
parameters('datastore my_url filter ctxsys.null_filter section group htmgroup');
```

See Also: ["Creating Preferences"](#) in this chapter for more examples on creating a custom context index.

Creating a CTXCAT Index

The CTXCAT indextype is well-suited for indexing small text fragments and related information. If created correctly, this type of index can give better structured query performance over a CONTEXT index.

CTXCAT Index and DML

A CTXCAT index is transactional. When you perform DML (inserts, updates, and deletes) on the base table, Oracle automatically synchronizes the index. Unlike a CONTEXT index, no CTX_DDL.SYNC_INDEX is necessary.

About CTXCAT Sub-Indexes and Their Costs

A CTXCAT index is comprised of sub-indexes that you define as part of your index set. You create a sub-index on one or more columns to improve mixed query performance.

However, adding sub-indexes to the index set has its costs. The time Oracle takes to create a CTXCAT index depends on its total size, and the total size of a CTXCAT index is directly related to

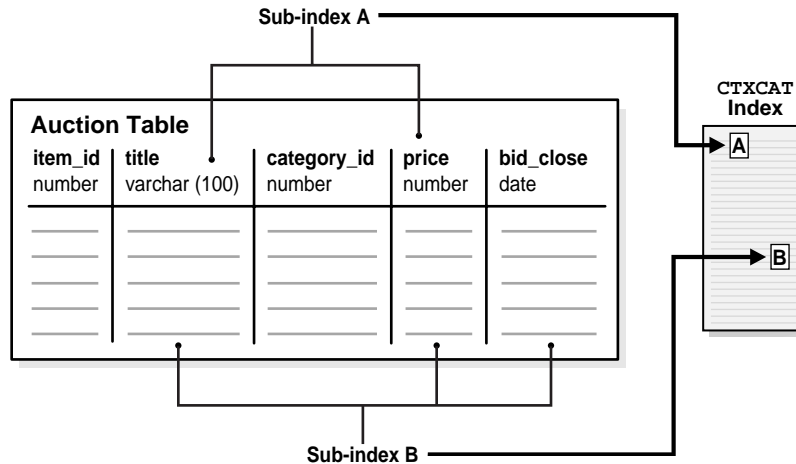
- total text to be indexed
- number of sub-indexes in the index set
- number of columns in the base table that make up the sub-indexes

Having many component indexes in your index set also degrades DML performance since more indexes must be updated.

Because of the added index time and disk space costs for creating a CTXCAT index, carefully consider the query performance benefit each component index gives your application before adding it to your index set.

Creating CTXCAT Sub-indexes

Figure 2–2



An online auction site that must store item descriptions, prices and bid-close dates for ordered look-up provides a good example for creating a CTXCAT index.

Figure 2–2 shows a table called AUCTION with the following schema:

```
create table auction(
  item_id number,
  title varchar2(100),
  category_id number,
  price number,
  bid_close date);
```

To create your sub-indexes, create an index set to contain them:

```
begin
  ctx_ddl.create_index_set('auction_iset');
end;
```

Next, determine the structured queries your application is likely to issue. The CATSEARCH query operator takes a mandatory text clause and optional structured clause.

In our example, this means all queries include a clause for the `title` column which is the text column.

Assume that the structured clauses fall into the following categories:

Structured Clauses	Sub-index Definition to Serve Query	Category
'price < 200'	'price'	A
'price = 150'		
'order by price'		
'price = 100 order by bid_close'	'price, bid_close'	B
'order by price, bid_close'		

Structured Query Clause Category A The structured query clause contains an expression for only the price column as follows:

```
SELECT FROM auction WHERE CATSEARCH(title, 'camera', 'price < 200') > 0;
SELECT FROM auction WHERE CATSEARCH(title, 'camera', 'price = 150') > 0;
SELECT FROM auction WHERE CATSEARCH(title, 'camera', 'order by price') > 0;
```

These queries can be served using sub-index B, but for efficiency you can also create a sub-index only on `price`, which we call sub-index A:

```
begin
  ctx_ddl.add_index('auction_iset', 'price'); /* sub-index A */
end;
```

Structured Query Clause Category B The structured query clause includes an equivalence expression for `price` ordered by `bid_close`, and an expression for ordering by `price` and `bid_close` in that order:

```
SELECT FROM auction WHERE CATSEARCH(title, 'camera', 'price = 100 order by bid_close') > 0;
SELECT FROM auction WHERE CATSEARCH(title, 'camera', 'order by price, bid_close') > 0;
```

These queries can be served with a sub-index defined as follows:

```
begin
  ctx_ddl.add_index('auction_iset', 'price, bid_close'); /* sub-index B */
end;
```

Like a combined b-tree index, the column order you specify with `CTX_DDL.ADD_INDEX` affects the efficiency and viability of the index scan Oracle uses to serve specific queries. For example, if two structured columns `p` and `q` have a b-tree index specified as `'p, q'`, Oracle cannot scan this index to sort `'order by q, p'`.

Creating CTXCAT Index

The following example combines the examples above and creates the index set preference with the three sub-indexes:

```
begin
  ctx_ddl.create_index_set('auction_iset');
  ctx_ddl.add_index('auction_iset', 'price'); /* sub-index A */
  ctx_ddl.add_index('auction_iset', 'price, bid_close'); /* sub-index B */
end;
```

[Figure 2-2](#) shows how the sub-indexes A and B are created from the auction table. Each sub-index is a b-tree index on the text column and the named structured columns. For example, sub-index A is an index on the `title` column and the `bid_close` column.

You create the combined catalog index with `CREATE INDEX` as follows:

```
CREATE INDEX auction_titlex ON AUCTION(title) INDEXTYPE IS CTXCAT PARAMETERS
('index set auction_iset');
```

See Also: *Oracle Text Reference* to learn more about creating a CTXCAT index with `CREATE INDEX`.

Creating a CTXRULE Index

You use the CTXRULE index to build a document classification application. You create a table of queries and then index them. With a CTXRULE index, you can use the `MATCHES` operator to classify single documents.

Create a Table of Queries

The first step is to create a table of queries that define your classifications. We create a table `myqueries` to hold the category name and query text:

```
CREATE TABLE myqueries (
  queryid NUMBER PRIMARY KEY,
  category VARCHAR2(30)
  query VARCHAR2(2000)
);
```

Populate the table with the classifications and the queries that define each. For example, consider a classification for the subjects *US Politics*, *Music*, and *Soccer*:

```
INSERT INTO myqueries VALUES(1, 'US Politics', 'democrat or republican');
INSERT INTO myqueries VALUES(2, 'Music', 'ABOUT(music)');
INSERT INTO myqueries VALUES(3, 'Soccer', 'ABOUT(soccer)');
```

Create the CTXRULE Index

Use CREATE INDEX to create the CTXRULE index. You can specify lexer, storage, section group, and wordlist parameters if needed:

```
CREATE INDEX ON myqueries(query) INDEXTYPE IS CTXRULE PARAMETERS('lexer lexer_
pref storage storage_pref section group section_pref wordlist wordlist_pref');
```

Note: The filter, memory, datastore, stoplist, and [no]populate parameters do not apply to the CTXRULE index type.

Classifying a Document

With a CTXRULE index created on query set, you can use the MATCHES operator to classify a document.

Assume that incoming documents are stored in the table news:

```
CREATE TABLE news (
    newsid NUMBER,
    author VARCHAR2(30),
    source VARCHAR2(30),
    article CLOB);
```

You can create a before insert trigger with MATCHES to route each document to another table news_route based on its classification:

```
BEGIN
    -- find matching queries
    FOR c1 IN (select category
              from myqueries
              where MATCHES(query, :new.article)>0)
    LOOP
        INSERT INTO news_route(newsid, category)
            VALUES (:new.newsid, c1.category);
    END LOOP;
END;
```

Index Maintenance

This section describes maintaining your index in the event of an error or indexing failure.

Viewing Index Errors

Sometimes an indexing operation might fail or not complete successfully. When the system encounters an error indexing a row, it logs the error in an Oracle Text view.

You can view errors on your indexes with `CTX_USER_INDEX_ERRORS`. View errors on all indexes as `CTXSYS` with `CTX_INDEX_ERRORS`.

For example to view the most recent errors on your indexes, you can issue:

```
SELECT err_timestamp, err_text FROM ctx_user_index_errors ORDER BY err_timestamp  
DESC;
```

To clear the view of errors, you can issue:

```
DELETE FROM ctx_user_index_errors;
```

See Also: *Oracle Text Reference* to learn more about these views.

Dropping an Index

You must drop an existing index before you can re-create it with `CREATE INDEX`.

You drop an index using the `DROP INDEX` command in SQL.

For example, to drop an index called `newsindex`, issue the following SQL command:

```
DROP INDEX newsindex;
```

If Oracle cannot determine the state of the index, for example as a result of an indexing crash, you cannot drop the index as described above. Instead use:

```
DROP INDEX newsindex FORCE;
```

See Also: *Oracle Text Reference* to learn more about this command.

Resuming Failed Index

You can resume a failed index creation operation using the `ALTER INDEX` command. You typically resume a failed index after you have investigated and corrected the index failure.

Index optimization commits at regular intervals. Therefore if an optimization operation fails, all optimization work has already been saved.

See Also: *Oracle Text Reference* to learn more about the ALTER INDEX command syntax.

Example: Resuming a Failed Index

The following command resumes the indexing operation on `newsindex` with 2 megabytes of memory:

```
ALTER INDEX newsindex REBUILD PARAMETERS('resume memory 2M');
```

Rebuilding an Index

You can rebuild a valid index using ALTER INDEX. You might rebuild an index when you want to index with a new preference.

See Also: *Oracle Text Reference* to learn more about the ALTER INDEX command syntax.

Example: Rebuilding and Index

The following command rebuilds the index, replacing the `lexer` preference with `my_lexer`.

```
ALTER INDEX newsindex REBUILD PARAMETERS('replace lexer my_lexer');
```

Dropping a Preference

You might drop a custom index preference when you no longer need it for indexing. You drop index preferences with the procedure `CTX_DDL.DROP_PREFERENCE`. Dropping a preference does not affect the index created from the preference.

See Also: *Oracle Text Reference* to learn more about the syntax for the `CTX_DDL.DROP_PREFERENCE` procedure.

Example

The following code drops the preference `my_lexer`.

```
begin
ctx_ddl.drop_preference('my_lexer');
end;
```

Managing DML Operations for a CONTEXT Index

DML operations to the base table refer to when documents are inserted, updated or deleted from the base table. This section describes how you can monitor, synchronize, and optimize the Oracle Text CONTEXT index when DML operations occur.

Note: CTXCAT indexes are transactional and thus updated immediately when there is an update to the base table. Manual synchronization as described in this section is not necessary for a CTXCAT index.

Viewing Pending DML

When documents in the base table are inserted, updated, or deleted, their ROWIDs are held in a DML queue until you synchronize the index. You can view this queue with the CTX_USER_PENDING view.

For example, to view pending DML on all your indexes, issue the following statement:

```
SELECT pnd_index_name, pnd_rowid, to_char(pnd_timestamp, 'dd-mon-yyyy
hh24:mi:ss') timestamp FROM ctx_user_pending;
```

This statement gives output in the form:

PND_INDEX_NAME	PND_ROWID	TIMESTAMP
MYINDEX	AAADXnAABAAAS3SAAC	06-oct-1999 15:56:50

See Also: *Oracle Text Reference* to learn more about this view.

Synchronizing the Index

Synchronizing the index involves processing all pending updates, inserts, and deletes to the base table. You can do this in PL/SQL with the CTX_DDL.SYNC_INDEX procedure.

The following example synchronizes the index with 2 megabytes of memory:

```
begin
  ctx_ddl.sync_index('myindex', '2M');
end;
```


Setting Background DML

You can set `CTX_DDL.SYNC_INDEX` to run automatically at regular intervals using the `DBMS_JOB.SUBMIT` procedure. Oracle Text includes a SQL script you can use to do this. The location of this script is:

```
$ORACLE_HOME/ctx/sample/script/drjobdml.sql
```

To use this script, you must be the index owner and you must have execute privileges on the `CTX_DDL` package. You must also set the `job_queue_processes` parameter in your Oracle initialization file.

For example, to set the index synchronization to run every 360 minutes on `myindex`, you can issue the following in SQL*Plus:

```
SQL> @drjobdml myindex 360
```

See Also: *Oracle Text Reference* to learn more about the `CTX_DDL.SYNC_INDEX` command syntax.

Index Optimization

Frequent index synchronization can fragment your CONTEXT index. Index fragmentation can adversely affect query response time. You can optimize your CONTEXT index to reduce fragmentation and index size and so improve query performance.

To understand index optimization, you must understand the structure of the index and what happens when it is synchronized.

CONTEXT Index Structure

The CONTEXT index is an inverted index where each word contains the list of documents that contain that word. For example, after a single initial indexing operation, the word `DOG` might have an entry as follows:

```
DOG DOC1 DOC3 DOC5
```

Index Fragmentation

When new documents are added to the base table, the index is synchronized by adding new rows. Thus if you add a new document (`DOC 7`) with the word *dog* to the base table and synchronize the index, you now have:

```
DOG DOC1 DOC3 DOC5
DOG DOC7
```

Subsequent DML will also create new rows:

```
DOG DOC1 DOC3 DOC5  
DOG DOC7  
DOG DOC9  
DOG DOC11
```

Adding new documents and synchronizing the index causes index fragmentation. In particular, background DML which synchronizes the index frequently generally produces more fragmentation than synchronizing in batch.

Less frequent batch processing results in longer document lists, reducing the number of rows in the index and hence reducing fragmentation.

You can reduce index fragmentation by optimizing the index in either FULL or FAST mode with `CTX_DDL.OPTIMIZE_INDEX`.

Document Invalidation and Garbage Collection

When documents are removed from the base table, Oracle Text marks the document as removed but does not immediately alter the index.

Because the old information takes up space and can cause extra overhead at query time, you must remove the old information from the index by optimizing it in FULL mode. This is called garbage collection.

Optimizing in FULL mode for garbage collection is necessary when you have frequent updates or deletes to the base table.

Single Token Optimization

In addition to optimizing the entire index, you can optimize single tokens. You can use token mode to optimize index tokens that are frequently searched, without spending time on optimizing tokens that are rarely referenced.

For example, you can specify that only the token *DOG* be optimized in the index, if you know that this token is updated and queried frequently.

An optimized token can improve query response time for the token.

To optimize an index in token mode, you can use `CTX_DDL.OPTIMIZE_INDEX`.

Examples: Optimizing the Index

To optimize an index, Oracle recommends that you use `CTX_DDL.OPTIMIZE_INDEX`.

See Also: *Oracle Text Reference* for the CTX_DDL.OPTIMIZE_INDEX command syntax and examples.

This chapter describes Oracle Text querying and associated features. The following topics are covered:

- [Overview of Queries](#)
- [Query Operators for CONTAINS](#)
- [Query Operators for CATSEARCH](#)
- [Optimizing for Response Time](#)
- [Counting Hits](#)

Overview of Queries

The basic Oracle Text query takes a query expression, usually a word with or without operators, as input. Oracle returns all documents (previously indexed) that satisfy the expression along with a relevance score for each document. Scores can be used to order the documents in the result set.

To issue an Oracle Text query, use the SQL `SELECT` statement with either the `CONTAINS` or `CATSEARCH` operator. You can use these operators programmatically wherever you can use the `SELECT` statement, such as in PL/SQL cursors.

Use the `MATCHES` operator to classify documents with a `CTXRULE` index.

Querying with `CONTAINS`

When you create an index of type `context`, you must use the `CONTAINS` operator to issue your query. An index of type `context` is suited for indexing collections of large coherent documents.

With the `CONTAINS` operator, you can use a number of operators to define your search criteria. These operators enable you to issue logical, proximity, fuzzy, stemming, thesaurus and wildcard searches. With a correctly configured index, you can also issue section searches on documents that have internal structure such as HTML and XML.

With `CONTAINS`, you can also use the `ABOUT` operator to search on document themes.

`CONTAINS` SQL Example

In the `SELECT` statement, specify the query in the `WHERE` clause with the `CONTAINS` operator. Also specify the `SCORE` operator to return the score of each hit in the hitlist. The following example shows how to issue a query:

```
SELECT SCORE(1) title from news WHERE CONTAINS(text, 'oracle', 1) > 0;
```

You can order the results from the highest scoring documents to the lowest scoring documents using the `ORDER BY` clause as follows:

```
SELECT SCORE(1), title from news
       WHERE CONTAINS(text, 'oracle', 1) > 0
       ORDER BY SCORE(1) DESC;
```

CONTAINS PL/SQL Example

In a PL/SQL application, you can use a cursor to fetch the results of the query.

The following example issues a `CONTAINS` query against the `NEWS` table to find all articles that contain the word *oracle*. The titles and scores of the first ten hits are output.

```
declare
  rowno number := 0;
begin
  for c1 in (SELECT SCORE(1) score, title FROM news
            WHERE CONTAINS(text, 'oracle', 1) > 0
            ORDER BY SCORE(1) DESC)
  loop
    rowno := rowno + 1;
    dbms_output.put_line(c1.title||': '||c1.score);
    exit when rowno = 10;
  end loop;
end;
```

This example uses a cursor `FOR` loop to retrieve the first ten hits. An alias *score* is declared for the return value of the `SCORE` operator. The score and title are output to standard out using cursor dot notation.

Structured Query with CONTAINS

A structured query, also called a mixed query, is a query that has a `CONTAINS` predicate to query a text column and has another predicate to query a structured data column.

To issue a structured query, you specify the structured clause in the `WHERE` condition of the `SELECT` statement.

For example, the following `SELECT` statement returns all articles that contain the word *oracle* that were written on or after October 1, 1997:

```
SELECT SCORE(1), title, issue_date from news
       WHERE CONTAINS(text, 'oracle', 1) > 0
       AND issue_date >= ('01-OCT-97')
       ORDER BY SCORE(1) DESC;
```

Note: Even though you can issue structured queries with CONTAINS, consider creating a `ctxcat` index and issuing the query with CATSEARCH, which offers better structured query performance.

Querying with CATSEARCH

When you create an index of type `ctxcat`, you must use the CATSEARCH operator to issue your query. An index of type `ctxcat` is best suited when your application stores short text fragments in the text column and other associated information in related columns.

For example, an application serving an online auction site might have a table that stores item description in a text column and associated information such as date and price in other columns. With a `ctxcat` index, you can create b-tree indexes on one or more of these columns. The result is that when you use the CATSEARCH operator to search a `ctxcat` index, query performance is generally faster for mixed queries.

The operators available for CATSEARCH queries are limited to logical operations such as AND or OR. The operators you can use to define your structured criteria are greater than, less than, equality, BETWEEN, and IN.

CATSEARCH SQL Query

A typical query with CATSEARCH might include a structured clause as follows to find all rows that contain the word *camera* ordered by the `bid_close` date:

```
SELECT FROM auction WHERE CATSEARCH(title, 'camera', 'order by bid_close desc')>
0;
```

The type of structured query you can issue depends on how you create your sub-indexes.

See Also: ["Creating a CTXCAT Index" in Chapter 2, "Indexing"](#).

CATSEARCH Structured Query

You specify the structured part of a CATSEARCH query with the `structured_query` parameter. The columns you name in the structured expression must have a corresponding sub-index.

For example, assuming that `category_id` and `bid_close` have a sub-index in the `ctxcat` index for the `AUCTION` table, you can issue the following structured query:

```
SELECT FROM auction WHERE CATSEARCH(title, 'camera', 'category_id=99 order by bid_close desc') > 0;
```

CATSEARCH PL/SQL Example

You can use a cursor to process the output of a CATSEARCH query as you do for `CONTAINS`.

Querying with MATCHES

When you create an index of type `CTXRULE`, you must use the `MATCHES` operator to classify your documents. The `CTXRULE` index is essentially an index on the set of queries that define your classifications.

For example, if you have an incoming stream of documents that need to be routed according to content, you can create a set of queries that define your categories. You create the queries as rows in a text column. You then index the table to create a `CTXRULE` index. When documents arrive, you use the `MATCHES` operator to classify each document.

MATCHES SQL Query

A `matches` query finds all rows in a query table that match a given document. Assuming that a table `querytable` has a `CTXRULE` index associated with it, you can issue the following query:

```
SELECT classification FROM querytable WHERE MATCHES(text, 'Smith is a common name in the United States') > 0;
```

MATCHES PL/SQL Example

The following example assumes that the table of queries `tdrbrn0101` has a `CTXRULE` index associated with it. It also assumes that the table `newsfeed` contains a set of news articles to be categorized.

This example loops through the `newsfeed` table, categorizing each article using the `MATCHES` operator. The results are stored in the `tdrbrn0102` table.

```
PROMPT Populate the category table based on newsfeed articles
PROMPT
set serveroutput on;
declare
  mypk      number;
  mytitle  varchar2(1000);
  myarticles clob;
  mycategory varchar2(100);
  cursor doccur is select pk,title,articles from newsfeed;
  cursor mycur is  select category from tdrbrn0101 where matches(rule,
myarticles)>0;
  cursor rescur is select category, pk, title from tdrbrn0102 order by
category,pk;

begin
  dbms_output.enable(1000000);
  open doccur;
  loop
    fetch doccur into mypk, mytitle, myarticles;
    exit when doccur%notfound;
    open mycur;
    loop
      fetch mycur into mycategory;
      exit when mycur%notfound;
      insert into tdrbrn0102 values(mycategory, mypk, mytitle);
    end loop;
    close mycur;
    commit;
  end loop;
  close doccur;
  commit;

end;
/
```

The following example displays the categorized articles by category.

```

PROMPT display the list of articles for every category
PROMPT

declare
  mypk number;
  mytitle varchar2(1000);
  mycategory varchar2(100);
  cursor catcur is select category from tdrbrn0101 order by category;
  cursor rescur is select pk, title from tdrbrn0102 where category=mycategory
order by pk;

begin
  dbms_output.enable(1000000);
  open catcur;
  loop
    fetch catcur into mycategory;
    exit when catcur%notfound;
    dbms_output.put_line('***** CATEGORY: '||mycategory||' *****');
  open rescur;
  loop
    fetch rescur into mypk, mytitle;
    exit when rescur%notfound;
  dbms_output.put_line('** ('||mypk||'). '||mytitle);
  end loop;
  close rescur;
  dbms_output.put_line('***');
  dbms_output.put_
line('*****');
  end loop;
  close catcur;
end;
/

```

Word and Phrase Queries

A word query is a query on a word or phrase. For example, to find all the rows in your text table that contain the word *dog*, you issue a query specifying *dog* as your query term.

You can issue word queries with both CONTAINS and CATSEARCH SQL operators.

If multiple words are contained in a query expression, separated only by blank spaces (no operators), the string of words is considered a phrase and Oracle searches for the entire string during a query.

For example, to find all documents that contain the phrase *international law*, you issue your query with the phrase *international law*.

Querying Stopwords

Stopwords are words for which Oracle does not create an index entry. They are usually common words in your language that are unlikely to be searched on by themselves.

Oracle Text includes a default list of stopwords for your language. This list is called a stoplist. For example, in English, the words *this* and *that* are defined as stopwords in the default stoplist. You can modify the default stoplist or create new stoplists with the CTX_DDL package. You can also add stopwords after indexing with the ALTER INDEX statement.

You cannot query on a stopword by itself or on a phrase composed of only stopwords. For example, a query on the word *this* returns no hits when *this* is defined as a stopword.

You can query on phrases that contain stopwords as well as non-stopwords such as *this boy talks to that girl*. This is possible because the Oracle Text index records the position of stopwords even though it does not create an index entry for them.

When you include a stopword within your query phrase, the stopword matches any word. For example, the query:

```
'Jack was big'
```

matches phrases such as *Jack is big* and *Jack grew big* assuming *was* is a stopword.

ABOUT Queries and Themes

An ABOUT query is a query on a document theme. A document theme is a concept that is sufficiently developed in the text. For example, an ABOUT query on *US politics* might return documents containing information about US presidential elections and US foreign policy. Documents need not contain the exact phrase *US politics* to be returned.

During indexing, document themes are derived from the knowledge base, which is a hierarchical list of categories and concepts that represents a view of the world. Some examples of themes in the knowledge catalog are concrete concepts such as *jazz music*, *football*, or *Nelson Mandela*. Themes can also be abstract concepts such as *happiness* or *honesty*.

During indexing, the system can also identify and index document themes that are sufficiently developed in the document, but do not exist in the knowledge base.

You can augment the knowledge base to define concepts and terms specific to your industry or query application. When you do so, ABOUT queries are more precise for the added concepts.

ABOUT queries perform best when you create a theme component in your index. Theme components are created by default for English and French.

Note: ABOUT queries are supported with only the CONTAINS operator.

Querying Stopthemes

Oracle enables you to query on themes with the ABOUT operator. A stoptheme is a theme that is not to be indexed. You can add and remove stopthemes with the CTX_DLL package. You can add stopthemes after indexing with the ALTER INDEX statement.

Query Expressions

A query expression is everything in between the single quotes in the `text_query` argument of the CONTAINS or CATSEARCH operator. What you can include in a query expression in a CONTAINS query is different from what you can include in a CATSEARCH operator.

CONTAINS Operator

A CONTAINS query expression can contain query operators that enable logical, proximity, thesaural, fuzzy, and wildcard searching. Querying with stored expressions is also possible. Within the query expression, you can use grouping characters to alter operator precedence.

With CONTAINS, you can also use the ABOUT query to query document themes.

See Also: ["Query Operators for CONTAINS"](#) in this chapter.

CATSEARCH Operator

With the CATSEARCH operator, you specify your query expression with the `text_query` operator and your optional structured criteria with the `structured_`

query argument. The `text_query` argument is limited to querying words and phrases. You can use logical operations, such as logical and, or, and not.

With `structured_query` argument, you specify your structured criteria. You can use the following SQL operations:

- =
- <=
- >=
- >
- <
- IN
- BETWEEN

You can also use ORDER BY clause to order your output.

See Also: ["Query Operators for CATSEARCH"](#) in this chapter.

MATCHES Operator

The MATCHES operator takes a document as input and finds all rows in a query table that match it. You do not specify query expressions in the MATCHES operator.

Case-Sensitive Searching

Oracle Text supports case-sensitivity for word and ABOUT queries.

Word Queries

Word queries are case-insensitive by default. This means that a query on the term *dog* returns the rows in your text table that contain the word *dog*, *Dog*, or *DOG*.

You can enable case-sensitive searching by enabling the `mixed_case` attribute in your BASIC_LEXER index preference. With a case-sensitive index, your queries must be issued in exact case. This means that a query on *Dog* matches only documents with *Dog*. Documents with *dog* or *DOG* are not returned as hits.

Stopwords and Case-Sensitivity If you have case-sensitivity enabled for word queries and you issue a query on a phrase containing stopwords and non-stopwords, you must specify the correct case for the stopwords. For example, a query on this boy

talks to that girl does not return text that contains the phrase *This boy talks to that girl*, assuming *this* is a stopword.

ABOUT Queries

ABOUT queries give the best results when your query is formulated with proper case. This is because the normalization of your query is based on the knowledge catalog which is case-sensitive. Attention to case is required especially for words that have different meanings depending on case, such as *turkey* the bird and *Turkey* the country.

However, you need not enter your query in exact case to obtain relevant results from an ABOUT query. The system does its best to interpret your query. For example, if you enter a query of ORACLE and the system does not find this concept in the knowledge catalog, the system might use Oracle as a related concept for look-up.

Query Feedback

Feedback information provides broader term, narrower term, and related term information for a specified query with a context index. You obtain this information programatically with the CTX_QUERY.HFEEDBACK procedure.

Broader term, narrower term, and related term information is useful for suggesting other query terms to the user in your query application.

The feedback information returned is obtained from the knowledge base and contains only those terms that are also in the index. This increases the chances that terms returned from HFEEDBACK produce hits over the currently indexed document set.

See Also: *Oracle Text Reference* for more information about using CTX_QUERY.HFEEDBACK

Query Explain Plan

Explain plan information provides a graphical representation of the parse tree for a CONTAINS query expression. You can obtain this information programatically with the CTX_QUERY.EXPLAIN procedure.

Explain plan information tells you how a query is expanded and parsed without having the system execute the query. Obtaining explain information is useful for knowing the expansion for a particular stem, wildcard, thesaurus, fuzzy, soundex, or ABOUT query. Parse trees also show the following information:

- order of execution
- ABOUT query normalization
- query expression optimization
- stop-word transformations
- breakdown of composite-word tokens for supported languages

See Also: *Oracle Text Reference* for more information about using CTX_QUERY.EXPLAIN

Query Operators for CONTAINS

With the CONTAINS operator, you can add complexity to your searches with query operators. You use the query operators in your query expression. For example, the logical operator AND allows you to search for all documents that contain two different words. The ABOUT operator allows you to search on concepts.

You can also use the WITHIN operator for section searching, the NEAR operator for proximity searches, the stem, fuzzy, and thesaural operators for expanding a query expression.

The following sections describe some of the Oracle Text operators.

See Also: *Oracle Text Reference* for complete information about using query operators.

ABOUT Query

Use the ABOUT operator in English or French to query on a concept. The query string is usually a concept or theme that represents the idea to be searched on. Oracle returns the documents that contain the theme.

Word information and theme information are combined into a single index. To issue a theme query, your index must have a theme component which is created by default in English and French.

You issue a theme query using the ABOUT operator inside the query expression. For example, to retrieve all documents that are about *politics*, write your query as follows:

```
SELECT SCORE(1), title FROM news
       WHERE CONTAINS(text, 'about(politics)', 1) > 0
       ORDER BY SCORE(1) DESC;
```

See Also: *Oracle Text Reference* for more information about using the ABOUT operator.

Logical Operators

Logical operators such as AND or OR allow you to limit your search criteria in a number of ways. The following table describes some of these operators.

Operator	Symbol	Description	Example Expression
AND	&	Use the AND operator to search for documents that contain at least one occurrence of <i>each</i> of the query terms. Score returned is the minimum of the operands.	'cats AND dogs' 'cats & dogs'
OR		Use the OR operator to search for documents that contain at least one occurrence of <i>any</i> of the query terms. Score returned is the maximum of the operands.	'cats dogs' 'cats OR dogs'
NOT	~	Use the NOT operator to search for documents that contain one query term and not another.	To obtain the documents that contain the term <i>animals</i> but not <i>dogs</i> , use the following expression: 'animals ~ dogs'
ACCUM	,	Use the ACCUM operator to search for documents that contain at least one occurrence of any of the query terms. The accumulate operator ranks documents according to the total term weight of a document.	The following query returns all documents that contain the terms <i>dogs</i> , <i>cats</i> and <i>puppies</i> giving the highest scores to the documents that contain all three terms: 'dogs, cats, puppies'
EQUIV	=	Use the EQUIV operator to specify an acceptable substitution for a word in a query.	The following example returns all documents that contain either the phrase <i>alsatians are big dogs</i> or <i>German shepherds are big dogs</i> : 'German shepherds=alsatians are big dogs'

Section Searching

Section searching is useful for when your document set is HTML or XML. For HTML, you can define sections using embedded tags and then use the WITHIN operator to search these sections.

For XML, you can have the system automatically create sections for you. You can query with the WITHIN operator or with the INPATH operator for path searching.

See Also: [Chapter 6, "Document Section Searching"](#)

Proximity Queries with NEAR Operator

You can search for terms that are near to one another in a document with the NEAR operator.

For example, to find all documents where *dog* is within 6 words of *cat*, issue the following query:

```
'near((dog, cat), 6)'
```

See Also: *Oracle Text Reference* for more information about using the NEAR operator.

Fuzzy, Stem, Soundex, Wildcard and Thesaurus Expansion Operators

You can expand your queries into longer word lists with operators such as wildcard, fuzzy, stem, soundex, and thesaurus.

See Also: *Oracle Text Reference* for more information about using these operators.

Stored Query Expressions

You can use the procedure CTX_QUERY.STORE_SQE to store the definition of a query without storing any results. Referencing the query with the CONTAINS SQE operator references the definition of the query. In this way, stored query expressions make it easy for defining long or frequently used query expressions.

Stored query expressions are not attached to an index. When you call CTX_QUERY.STORE_SQE, you specify only the name of the stored query expression and the query expression.

The query definitions are stored in the Text data dictionary. Any user can reference a stored query expression.

See Also: *Oracle Text Reference* to learn more about the syntax of `CTX_QUERY.STORE_SQE`.

Defining a Stored Query Expression

You define and use a stored query expression as follows:

1. Call `CTX_QUERY.STORE_SQE` to store the results for the text column. With `STORE_SQE`, you specify a name for the stored query expression and a query expression.
2. Call the stored query expression in a query expression using the `SQE` operator. Oracle returns the results of the stored query expression in the same way it returns the results of a regular query. The query is evaluated at the time the stored query expression is called.

You can delete a stored query expression using `REMOVE_SQE`.

SQE Example

The following example creates a stored query expression called *disaster* that searches for documents containing the words *tornado*, *hurricane*, or *earthquake*:

```
begin
ctx_query.store_sqe('disaster', 'tornado | hurricane | earthquake');
end;
```

To execute this query in an expression, write your query as follows:

```
SELECT SCORE(1), title from news
WHERE CONTAINS(text, 'SQE(disaster)', 1) > 0
ORDER BY SCORE(1);
```

See Also: *Oracle Text Reference* to learn more about the syntax of `CTX_QUERY.STORE_SQE`.

Calling PL/SQL Functions in CONTAINS

You can call user-defined functions directly in the `CONTAINS` clause as long as the function satisfies the requirements for being named in a SQL statement. The caller must also have `EXECUTE` privilege on the function.

For example, assuming the function *french* returns the French equivalent of an English word, you can search on the French word for *cat* by writing:

```
SELECT SCORE(1), title from news
WHERE CONTAINS(text, french('cat'), 1) > 0
```

```
ORDER BY SCORE(1);
```

See Also: *Oracle9i SQL Reference* for more information about creating user functions and calling user functions from SQL,

Query Operators for CATSEARCH

The CATSEARCH operator has a simpler query language than CONTAINS. Its query language supports logical operations such as AND and OR as well as phrase queries.

The CATSEARCH query operators have the following syntax:

Operation	Syntax	Description of Operation
Logical AND	a b c	Returns rows that contain a, b and c.
Logical OR	a b c	Returns rows that contain a, b, or c.
Logical NOT	a - b	Returns rows that contain a and not b.
hyphen with no space	a-b	Hyphen treated as a regular character. For example, if the hyphen is defined as skipjoin, words such as <i>web-site</i> treated as the single query term <i>website</i> . Likewise, if the hyphen is defined as a printjoin, words such as <i>web-site</i> treated as <i>web site</i> with the space in the CTXCAT query language.
" "	"a b c"	Returns rows that contain the phrase "a b c". For example, entering "Sony CD Player" means return all rows that contain this sequence of words.
()	(A B) C	Parentheses group operations. This query is equivalent to the CONTAINS query (A &B) C.

See Also: *Oracle Text Reference* for more information about using these operators.

Optimizing for Response Time

A query optimized for response time provides a fast solution for when you need the highest scoring documents from a hitlist.

The example below returns the first twenty hits to standard out. This example uses the `FIRST_ROWS` hint and a cursor.

```
declare
cursor c is
  select /*+ FIRST_ROWS */ title, score(1) score
    from news
   where contains(txt_col, 'dog', 1) > 0
   order by score(1) desc;
begin
  for c1 in c
  loop
    dbms_output.put_line(c1.score||':'||substr(c1.title,1,50));
    exit when c%rowcount = 20;
  end loop;
end;
/
```

See Also: [Chapter 5, "Query Tuning"](#)

Retrieving a Range of Documents

A query optimized for response time provides a fast solution for when you need a range of documents from a hitlist sorted by score.

The solution uses the `FIRST_ROWS` hint in a cursor. The code loops through the cursor to process only the hits in the required range. The example below returns the sorted documents 11 to 20.

```
declare
cursor c is
  select /*+ FIRST_ROWS */ title, score(1) score
    from news
   where contains(txt_col, 'dog', 1) > 0
   order by score(1) desc;
begin
  for c1 in c
  loop
    if (c%rowcount > 10) then
      dbms_output.put_line(c1.score||':'||substr(c1.title,1,50));
    end if;
  end loop;
end;
```

```
        end if;  
        exit when c%rowcount = 20;  
    end loop;  
end;  
/
```

See Also: [Chapter 5, "Query Tuning"](#)

Counting Hits

To count the number of hits returned from a query with only a CONTAINS predicate, you can use CTX_QUERY.COUNT_HITS in PL/SQL or COUNT(*) in a SQL SELECT statement.

If you want a rough hit count, you can use CTX_QUERY.COUNT_HITS in estimate mode (EXACT parameter set to FALSE). With respect to response time, this is the fastest count you can get.

To count the number of hits returned from a query that contains a structured predicate, use the COUNT(*) function in a SELECT statement.

SQL Count Hits Example

To find the number of documents that contain the word *oracle*, issue the query with the SQL COUNT function as follows:

```
SELECT count(*) FROM news WHERE CONTAINS(text, 'oracle', 1) > 0;
```

Counting Hits with a Structured Predicate

To find the number of documents returned by a query with a structured predicate, use COUNT(*) as follows:

```
SELECT COUNT(*) FROM news WHERE CONTAINS(text, 'oracle', 1) > 0 and author = 'jones';
```

PL/SQL Count Hits Example

To find the number of documents that contain the word *oracle*, use COUNT_HITS as follows:

```
declare count number;
begin
  count := ctx_query.count_hits(index_name=>my_index, text_query=>'oracle',
                               exact => TRUE);
  dbms_output.put_line('Number of docs with oracle:');
  dbms_output.put_line(count);
end;
```

See Also: *Oracle Text Reference* to learn more about the syntax of CTX_QUERY.COUNT_HITS.

Document Presentation

This chapter describes document presentation. The following topics are covered:

- [Highlighting Query Terms](#)
- [Obtaining List of Themes, Gists, and Theme Summaries](#)

Highlighting Query Terms

In Oracle Text query applications, you can present selected documents with query terms highlighted for text queries or with themes highlighted for ABOUT queries.

You can generate three types of output associated with highlighting: a marked-up version of the document, a plain text version of the document (filtered output), and highlight offset information for the document.

The three types of output are generated by three different procedures in the CTX_DOC (document services) PL/SQL package. In addition, you can obtain plain text and HTML versions for each type of output.

Text highlighting

For text highlighting, you supply the query, and Oracle highlights words in document that satisfy the query. You can obtain plain-text or HTML highlighting.

Theme Highlighting

For ABOUT queries, the CTX_DOC procedures highlight and mark up words or phrases that best represent the ABOUT query.

CTX_DOC Highlighting Procedures

There are three highlighting procedures in CTX_DOC:

- HIGHLIGHT
- MARKUP
- FILTER

Highlight Procedure

Highlight offset information is useful for when you write your own custom routines for displaying documents.

To obtain highlight offset information, use the CTX_DOC.HIGHLIGHT procedure. This procedure takes a query and a document, and returns highlight offset information for either plaintext or HTML formats.

With offset information, you can display a highlighted version of document as desired. For example, you can display the document with different font types or colors rather than using the standard plain text markup obtained from CTX_DOC.MARKUP.

See Also: *Oracle Text Reference* for more information about using CTX_DOC.HIGHLIGHT.

Markup Procedure

The CTX_DOC.MARKUP procedure takes a document reference and a query, and returns a marked-up version of the document. The output can be either marked-up plaintext or marked-up HTML.

You can customize the markup sequence for HTML navigation.

See Also: *Oracle Text Reference* for more information about CTX_DOC.MARKUP.

Filter Procedure

When documents are stored in their native formats such as Microsoft Word, you can use the filter procedure CTX_DOC.FILTER to obtain either a plain text or HTML version of the document.

See Also: *Oracle Text Reference* for more information about CTX_DOC.FILTER.

Obtaining List of Themes, Gists, and Theme Summaries

The following table describes list of themes, gists, and theme summaries.

Table 4–1

Output Type	Description
List of Themes	A list of the main concepts of a document. You can generate list of themes where each theme is a single word or phrase or where each theme is a hierarchical list of parent themes.
Gist	Text in a document that best represents what the document is about as a whole.
Theme Summary	Text in a document that best represents a given theme in the document.

To obtain this output, you use procedures in the CTX_DOC supplied package. With this package, you can do the following:

- Identify documents by ROWID in addition to primary key
- Store results in-memory for improved performance

List of Themes

A list of themes is a list of the main concepts in a document. Use the CTX_DOC.THEMES procedure to generate lists of themes.

See Also: *Oracle Text Reference* to learn more about the command syntax for CTX_DOC.THEMES.

In-Memory Themes

The following example generates the top 10 themes for document 1 and stores them in an in-memory table called `the_themes`. The example then loops through the table to display the document themes.

```
declare
  the_themes ctx_doc.theme_tab;

begin
  ctx_doc.themes('myindex','1',the_themes, numthemes=>10);
  for i in 1..the_themes.count loop
    dbms_output.put_line(the_themes(i).theme||':'||the_themes(i).weight);
  end loop;
```

```
end;
```

Result Table Themes

To create a theme table:

```
create table ctx_themes (query_id number,  
                        theme varchar2(2000),  
                        weight number);
```

Single Themes To obtain a list of themes where each element in the list is a single theme, issue:

```
begin  
ctx_doc.themes('newsindex',34,'CTX_THEMES',1,full_themes => FALSE);  
end;
```

Full Themes To obtain a list of themes where each element in the list is a hierarchical list of parent themes, issue:

```
begin  
ctx_doc.themes('newsindex',34,'CTX_THEMES',1,full_themes => TRUE);  
end;
```

Gist and Theme Summary

A gist is the text of a document that best represents what the document is about as a whole. A theme summary is the text of a document that best represents a single theme in the document.

Use the procedure `CTX_DOC.GIST` to generate gists and theme summaries. You can specify the size of the gist or theme summary when you call the procedure.

See Also: *Oracle Text Reference* to learn about the command syntax for `CTX_DOC.GIST`.

In-Memory Gist

The following example generates a non-default size generic gist of at most 10 paragraphs. The result is stored in memory in a CLOB locator. The code then de-allocates the returned CLOB locator after using it.

```
declare  
  gklob clob;  
  amt number := 40;
```

```
line varchar2(80);

begin
  ctx_doc.gist('newsindex', '34', 'gklob', 1, glevel => 'P', pov => 'GENERIC',
numParagraphs => 10);
  -- gklob is NULL when passed-in, so ctx-doc.gist will allocate a temporary
  -- CLOB for us and place the results there.

  dbms_lob.read(gklob, amt, 1, line);
  dbms_output.put_line('FIRST 40 CHARS ARE: ' || line);
  -- have to de-allocate the temp lob
  dbms_lob.freetemporary(gklob);
end;
```

Result Table Gists

To create a gist table:

```
create table ctx_gist (query_id number,
                      pov      varchar2(80),
                      gist      CLOB);
```

The following example returns a default sized paragraph level gist for document 34:

```
begin
ctx_doc.gist('newsindex', 34, 'CTX_GIST', 1, 'PARAGRAPH', pov => 'GENERIC');
end;
```

The following example generates a non-default size gist of ten paragraphs:

```
begin
ctx_doc.gist('newsindex', 34, 'CTX_GIST', 1, 'PARAGRAPH', pov => 'GENERIC',
numParagraphs => 10);
end;
```

The following example generates a gist whose number of paragraphs is ten percent of the total paragraphs in document:

```
begin
ctx_doc.gist('newsindex', 34, 'CTX_GIST', 1, 'PARAGRAPH', pov => 'GENERIC',
maxPercent => 10);
end;
```

Theme Summary

The following example returns a theme summary on the theme of *insects* for document with textkey 34. The default Gist size is returned.


```
begin
ctx_doc.gist('newsindex',34,'CTX_GIST',1, 'PARAGRAPH', pov => 'insects');
end;
```


5

Query Tuning

This appendix discusses how to tune your queries for better response time. The following topics are covered:

- [Optimizing Queries with Statistics](#)
- [Optimizing Queries for Response Time](#)
- [Optimizing Queries for Throughput](#)
- [Tuning Queries with Blocking Operations](#)

Optimizing Queries with Statistics

Query optimization with statistics uses the collected statistics on the tables and indexes in a query to select an execution plan that can process the query in the most efficient manner. The optimizer attempts to choose the best execution plan based on the following parameters:

- the selectivity on the CONTAINS predicate
- the selectivity of other predicates in the query
- the CPU and I/O costs of processing the CONTAINS predicates

The following sections describe how to use statistics with the extensible query optimizer. Optimizing with statistics allows for a more accurate estimation of the selectivity and costs of the CONTAINS predicate and thus a better execution plan.

Collecting Statistics

By default, Oracle uses the cost-based optimizer to determine the best execution plan for a query. To allow the optimizer to better estimate costs, you can calculate the statistics on the table you query. To do so, issue the following statement:

```
ANALYZE TABLE <table_name> COMPUTE STATISTICS;
```

Alternatively, you can estimate the statistics on a sample of the table as follows:

```
ANALYZE TABLE <table_name> ESTIMATE STATISTICS 1000 ROWS;
```

or

```
ANALYZE TABLE <table_name> ESTIMATE STATISTICS 50 PERCENT;
```

These statements collect statistics on all the objects associated with `table_name` including the table columns and any indexes (b-tree, bitmap, or Text domain) associated with the table. To re-collect the statistics on a table, you can issue the ANALYZE command as many times as necessary or use the DBMS_STATS package

See Also: *Oracle9i SQL Reference* and *Oracle9i Database Performance Guide and Reference* for more information about the ANALYZE command.

By collecting statistics on the Text domain index, the Oracle cost-based optimizer is able to do the following:

- estimate the selectivity of the CONTAINS predicate

- estimate the I/O and CPU costs of using the Text index, that is, the cost of processing the CONTAINS predicate using the domain index
- estimate the I/O and CPU costs of each invocation of CONTAINS

Knowing the selectivity of a CONTAINS predicate is useful for queries that contain more than one predicate, such as in structured queries. This way the cost-based optimizer can better decide whether to use the domain index to evaluate CONTAINS or to apply the CONTAINS predicate as a post filter.

Example

Consider the following structured query:

```
select score(1) from tab where contains(txt, 'freedom', 1) > 0 and author =  
'King' and year > 1960;
```

Assume the `author` column is of type `VARCHAR2` and the `year` column is of type `NUMBER`. Assume that there is a b-tree index on the `author` column.

Also assume that the structured `author` predicate is highly selective with respect to the CONTAINS predicate and the year predicate. That is, the structured predicate (`author = 'King'`) returns a much smaller number of rows with respect to the year and CONTAINS predicates individually, say 5 rows versus 1000 and 1500 rows respectively.

In this situation, Oracle can execute this query more efficiently by first doing a b-tree index range scan on the structured predicate (`author = 'King'`), followed by a table access by rowid, and then applying the other two predicates to the rows returned from the b-tree table access.

Note: When statistics are not collected for a Text index, the cost-based optimizer assumes low selectivity and index costs for the CONTAINS predicate.

Re-Collecting Statistics

After synchronizing your index, you can re-collect statistics on a single index to update the cost estimates. To do so, you can issue any of the following statements:

```
ANALYZE INDEX <index_name> COMPUTE STATISTICS;  
or
```

```
ANALYZE INDEX <index_name> ESTIMATE STATISTICS SAMPLE 1000 ROWS;
```

or

```
ANALYZE INDEX <index_name> ESTIMATE STATISTICS SAMPLE 50 PERCENT;
```

Deleting Statistics

You can delete the statistics associated with a table by issuing:

```
ANALYZE TABLE <table_name> DELETE STATISTICS;
```

You can delete statistics on one index by issuing the following statement:

```
ANALYZE INDEX <index_name> DELETE STATISTICS;
```

Optimizing Queries for Response Time

By default, Oracle optimizes queries for throughput. This results in queries returning all rows in shortest time possible.

However, in many cases, especially in a web-application scenario, queries must be optimized for response time, when you are only interested in obtaining the first few hits of a potentially large hitlist in the shortest time possible.

The following sections describe how to optimize Text queries for response time. You can do so in two ways:

- using `FIRST_ROWS` hint
- using `CHOOSE` and `DOMAIN_INDEX_SORT` hints

Note: Although both methods optimize for response time, the execution plans of the two methods obtained with `EXPLAIN PLAN` might be different for a given query.

Better Response Time with `FIRST_ROWS`

You can change the default query optimizer mode to optimize for response time using the `FIRST_ROWS` hint. When queries are optimized for response time, Oracle returns the first rows in the shortest time possible.

For example, consider the following PL/SQL block that uses a cursor to retrieve the first 20 hits of a query and uses the `FIRST_ROWS` hint to optimize the response time:

```
declare
cursor c is
select /*+ FIRST_ROWS */ pk, score(1), col from ctx_tab
      where contains(txt_col, 'test', 1) > 0 order by score(1) desc;
begin
for i in c
loop
insert into t_s values(i.pk, i.col);
exit when c%rowcount > 21;
end loop;
end;
/
```

The cursor `c` is a `SELECT` statement that returns the rowids that contain the word *test* in sorted order. The code loops through the cursor to extract the first 20 rows. These rows are stored in the temporary table `t_s`.

With the `FIRST_ROWS` hint, Oracle instructs the Text index to return rowids in score-sorted order, if possible.

Without the hint, Oracle sorts the rowids after the Text index has returned *all* the rows in unsorted order that satisfy the `CONTAINS` predicate. Retrieving the entire result set as such takes time.

Since only the first 20 hits are needed in this query, using the hint results in better performance.

Note: Use the `FIRST_ROWS` hint when you need only the first few hits of a query. When you need the entire result set, do not use this hint as it might result in poor performance.

Other Behavior with `FIRST_ROWS`

Besides instructing the Text index to return hits in score-sorted order, the `FIRST_ROWS` hint also tries to avoid blocking operations when optimizing queries for response time. Blocking operations include merge joins, hash joins, and bitmap operations.

As a result, using the `FIRST_ROWS` hint to optimize for response time might result in a different execution plan than using `CHOOSE` with `DOMAIN_INDEX_SORT`, which also optimizes for response time.

You can examine query execution plans using the `EXPLAIN PLAN` command in `SQL`.

See Also: *Oracle9i Database Performance Guide and Reference* for more information about the query optimizer and using hints such as `FIRST_ROWS` and `CHOOSE`.

For more information about the `EXPLAIN PLAN` command, *Oracle9i Database Performance Guide and Reference* and *Oracle9i SQL Reference*.

Better Response Time with `CHOOSE`

When you use the `CHOOSE` or `ALL_ROWS` optimizer hints, the query is optimized for throughput. This is the default optimizer mode. In this mode, Oracle does not

instruct the Text domain index to return score-sorted rows, choosing instead to sort all the rows fetched from the Text index.

To optimize for fast response time under CHOOSE or ALL_ROWS modes, you can use the DOMAIN_INDEX_SORT hint as follows:

```
declare
cursor c is
select /*+ CHOOSE DOMAIN_INDEX_SORT */ pk, score(1), col from ctx_tab
      where contains(txt_col, 'test', 1) > 0 order by score(1) desc;
begin
for i in c
loop
insert into t_s values(i.pk, i.col);
exit when c%rowcount > 21;
end loop;
end;
/
```

Note: Although you can optimize for response time with this method as well as with FIRST_ROWS by itself, the actual execution plans of the two methods obtained with EXPLAIN PLAN might be different for a given query.

See Also: *Oracle9i Database Performance Guide and Reference* for more information about the query optimizer and using hints such as FIRST_ROWS and CHOOSE.

For more information about the EXPLAIN PLAN command, *Oracle9i Database Performance Guide and Reference* and *Oracle9i SQL Reference*.

Optimizing Queries for Throughput

Optimizing a query for throughput returns all hits in the shortest time possible. This is the default behavior.

The following sections describe how you can explicitly optimize for throughput.

CHOOSE and ALL ROWS Modes

By default, queries are optimized for throughput under the CHOOSE and ALL_ROWS modes. When queries are optimized for throughput, Oracle returns *all* rows in the shortest time possible.

FIRST_ROWS Mode

In FIRST_ROWS mode, the Oracle optimizer optimizes for fast response time by having the Text domain index return score-sorted rows, if possible. This is the default behavior when you use the FIRST_ROWS hint.

If you want to optimize for better throughput under FIRST_ROWS, you can use the DOMAIN_INDEX_NO_SORT hint. Better throughput means you are interested in getting all the rows to a query in the shortest time.

The following example achieves better throughput by not using the Text domain index to return score-sorted rows. Instead, Oracle sorts the rows after all the rows that satisfy the CONTAINS predicate are retrieved from the index:

```
select /*+ FIRST_ROWS DOMAIN_INDEX_NO_SORT */ pk, score(1), col from ctx_tab
       where contains(txt_col, 'test', 1) > 0 order by score(1) desc;
```

See Also: *Oracle9i Database Performance Guide and Reference* for more information about the query optimizer and using hints such as FIRST_ROWS and CHOOSE.

Tuning Queries with Blocking Operations

Issuing a query with more than one predicate can cause a blocking operation in the execution plan. For example, consider the following mixed query:

```
select docid from myindex where contains(text, 'oracle', 1) > 0
AND colA > 5
AND colB > 1
AND colC > 3;
```

Assume that all predicates are unselective and colA, colB, and colC have bitmap indexes. The Oracle cost-based optimizer chooses the following execution plan:

```
TABLE ACCESS BY ROWIDS
  BITMAP CONVERSION TO ROWIDS
    BITMAP AND
      BITMAP INDEX COLA_BMX
      BITMAP INDEX COLB_BMX
      BITMAP INDEX COLC_BMX
    BITMAP CONVERSION FROM ROWIDS
      SORT ORDER BY
        DOMAIN INDEX MYINDEX
```

Since the BITMAP AND is a blocking operation, Oracle must temporarily save the rowid and score pairs returned from the Oracle Text domain index before executing the BITMAP AND operation.

Oracle attempts to save these rowid and score pairs in memory. However, when the size of the result set containing these rowid and score pairs exceeds the SORT_AREA_SIZE initialization parameter, Oracle spills these results to temporary segments on disk.

Since saving results to disk causes extra overhead, you can improve performance by increasing the SORT_AREA_SIZE parameter using ALTER SESSION as follows:

```
alter session set SORT_AREA_SIZE = <new memory size in bytes>;
```

For example, to set the buffer to approximately 8 megabytes, you can issue:

```
alter session set SORT_AREA_SIZE = 8300000;
```

See Also: *Oracle9i Database Performance Guide and Reference* and *Oracle9i Database Reference* for more information on SORT_AREA_SIZE.

Document Section Searching

This chapter describes how to use document sections in an Oracle Text query application.

The following topics are discussed in this chapter:

- [About Document Section Searching](#)
- [HTML Section Searching](#)
- [XML Section Searching](#)

About Document Section Searching

Section searching enables you to narrow text queries down to blocks of text within documents. Section searching is useful when your documents have internal structure, such as HTML and XML documents.

You can also search for text at the sentence and paragraph level.

Enabling Section Searching

The steps for enabling section searching for your document collection are:

1. Create a section group
2. Define your sections
3. Index your documents
4. Section search with WITHIN, INPATH, or HASPATH operators

Create a Section Group

Section searching is enabled by defining section groups. You use one of the system-defined section groups to create an instance of a section group. Choose a section group appropriate for your document collection.

You use section groups to specify the type of document set you have and implicitly indicate the tag structure. For instance, to index HTML tagged documents, you use the HTML_SECTION_GROUP. Likewise, to index XML tagged documents, you can use the XML_SECTION_GROUP.

The following table lists the different types of section groups you can use:

Section Group Preference	Description
NULL_SECTION_GROUP	This is the default. Use this group type when you define no sections or when you define <i>only</i> SENTENCE or PARAGRAPH sections.
BASIC_SECTION_GROUP	Use this group type for defining sections where the start and end tags are of the form <A> and .
HTML_SECTION_GROUP	Use this group type for indexing HTML documents and for defining sections in HTML documents.
XML_SECTION_GROUP	Use this group type for indexing XML documents and for defining sections in XML documents.

Section Group Preference	Description
AUTO_SECTION_GROUP	<p>Use this group type to automatically create a zone section for each start-tag/end-tag pair in an XML document. The section names derived from XML tags are case-sensitive as in XML.</p> <p>Attribute sections are created automatically for XML tags that have attributes. Attribute sections are named in the form attribute@tag.</p> <p>Stop sections, empty tags, processing instructions, and comments are not indexed.</p> <p>The following limitations apply to automatic section groups:</p> <ul style="list-style-type: none"> ■ You cannot add zone, field or special sections to an automatic section group. ■ Automatic sectioning does not index XML document types (root elements.) However, you can define stop-sections with document type. ■ The length of the indexed tags including prefix and namespace cannot exceed 64 characters. Tags longer than this are not indexed.
PATH_SECTION_GROUP	<p>Use this group type to index XML documents. Behaves like the AUTO_SECTION_GROUP.</p> <p>The difference is that with this section group you can do path searching with the INPATH and HASPATH operators. Queries are also case-sensitive for tag and attribute names.</p>
NEWS_SECTION_GROUP	<p>Use this group for defining sections in newsgroup formatted documents according to RFC 1036.</p>

You use the CTX_DDL package to create section groups and define sections as part of section groups. For example, to index HTML documents, create a section group with HTML_SECTION_GROUP:

```
begin
ctx_ddl.create_section_group('htmgroup', 'HTML_SECTION_GROUP');
end;
```

Define Your Sections

You define sections as part of the section group. The following example defines an zone section called heading for all text within the HTML < H1> tag:

```
begin
ctx_ddl.create_section_group('htmgroup', 'HTML_SECTION_GROUP');
ctx_ddl.add_zone_section('htmgroup', 'heading', 'H1');
end;
```

Note: If you are using the AUTO_SECTION_GROUP or PATH_SECTION_GROUP to index an XML document collection, you need not explicitly define sections since the system does this for you during indexing.

See Also: ["Section Types"](#) in this chapter for more information about sections.

["XML Section Searching"](#) in this chapter for more information about section searching with XML.

Index your Documents

When you index your documents, you specify your section group in the parameter clause of CREATE INDEX.

```
create index myindex on docs(htmlfile) indextype is ctxsys.context
parameters('filter ctxsys.null_filter section group htmgroup');
```

Section Searching with WITHIN Operator

When your documents are indexed, you can query within sections using the WITHIN operator. For example, to find all the documents that contain the word Oracle within their headings, issue the following query:

```
'Oracle WITHIN heading'
```

See Also: *Oracle Text Reference* to learn more about using the WITHIN operator.

Path Searching with INPATH and HASPATH Operators

When you use the PATH_SECTION_GROUP, the system automatically creates XML sections for you. In addition to using the WITHIN operator to issue queries, you can issue path queries with the INPATH and HASPATH operators.

See Also: ["XML Section Searching"](#) to learn more about using these operators.

Oracle Text Reference to learn more about using the INPATH operator.

Section Types

All sections types are blocks of text in a document. However, sections can differ in the way they are delimited and the way they are recorded in the index. Sections can be one of the following:

- zone section
- field section
- attribute section (for XML documents)
- special (sentence or paragraphs)

Zone Section

A zone section is a body of text delimited by start and end tags in a document. The positions of the start and end tags are recorded in the index so that any words in between the tags are considered to be within the section. Any instance of a zone section must have a start and an end tag.

For example, the text between the <TITLE> and </TITLE> tags can be defined as a zone section as follows:

```
<TITLE>Tale of Two Cities</TITLE>  
It was the best of times...
```

Zone sections can nest, overlap, and repeat within a document.

When querying zone sections, you use the WITHIN operator to search for a term across all sections. Oracle returns those documents that contain the term within the defined section.

Zone sections are well suited for defining sections in HTML and XML documents. To define a zone section, use CTX_DDL.ADD_ZONE_SECTION.

For example, assume you define the section *booktitle* as follows:

```
begin
ctx_ddl.create_section_group('htmgroup', 'HTML_SECTION_GROUP');
ctx_ddl.add_zone_section('htmgroup', 'booktitle', 'TITLE');
end;
```

After you index, you can search for all the documents that contain the term *Cities* within the section *booktitle* as follows:

```
'Cities WITHIN booktitle'
```

With multiple query terms such as (*dog and cat*) *WITHIN booktitle*, Oracle returns those documents that contain *cat* and *dog* within the same instance of a *booktitle* section.

Repeated Zone Sections Zone sections can repeat. Each occurrence is treated as a separate section. For example, if `<H1>` denotes a heading section, they can repeat in the same documents as follows:

```
<H1> The Brown Fox </H1>
<H1> The Gray Wolf </H1>
```

Assuming that these zone sections are named *Heading*, the query *Brown WITHIN Heading* returns this document. However, a query of (*Brown and Gray*) *WITHIN Heading* does not.

Overlapping Zone Sections Zone sections can overlap each other. For example, if `` and `<I>` denote two different zone sections, they can overlap in a document as follows:

```
plain <B> bold <I> bold and italic </B> only italic </I> plain
```

Nested Zone Sections Zone sections can nest, including themselves as follows:

```
<TD> <TABLE><TD>nested cell</TD></TABLE></TD>
```

Using the *WITHIN* operator, you can write queries to search for text in sections within sections. For example, assume the *BOOK1*, *BOOK2*, and *AUTHOR* zone sections occur as follows in documents *doc1* and *doc2*:

doc1:

```
<book1> <author>Scott Tiger</author> This is a cool book to read.</book1>
```

doc2:

```
<book2> <author>Scott Tiger</author> This is a great book to read.<book2>
```

Consider the nested query:

```
'Scott within author within book1'
```

This query returns only doc1.

Field Section

A field section is similar to a zone section in that it is a region of text delimited by start and end tags. A field section is different from a zone section in that the region is indexed separate from the rest of the document.

Since field sections are indexed differently, you can also get better query performance over zone sections for when you have a large number of documents indexed.

Field sections are more suited to when you have a single occurrence of a section in a document such as a field in a news header. Field sections can also be made visible to the rest of the document.

Unlike zone sections, field sections have the following restrictions:

- field sections cannot overlap
- field sections cannot repeat
- field sections cannot nest

Visible and Invisible Field Sections By default, field sections are indexed as a sub-document separate from the rest of the document. As such, field sections are invisible to the surrounding text and can only be queried by explicitly naming the section in the WITHIN clause.

You can make field sections visible if you want the text within the field section to be indexed as part of the enclosing document. Text within a visible field section can be queried with or without the WITHIN operator.

The following example shows the difference between using invisible and visible field sections.

The following code defines a section group `basicgroup` of the `BASIC_SECTION_GROUP` type. It then creates a field section in `basicgroup` called `Author` for the `<A>` tag. It also sets the visible flag to `FALSE` to create an invisible section:

```
begin
ctx_ddl_create_section_group('basicgroup', 'BASIC_SECTION_GROUP');
ctx_ddl.add_field_section('basicgroup', 'Author', 'A', FALSE);
end;
```

Because the `Author` field section is not visible, to find text within the `Author` section, you must use the `WITHIN` operator as follows:

```
'(Martin Luther King) WITHIN Author'
```

A query of *Martin Luther King* without the `WITHIN` operator does not return instances of this term in field sections. If you want to query text within field sections without specifying `WITHIN`, you must set the visible flag to `TRUE` when you create the section as follows:

```
begin
ctx_ddl.add_field_section('basicgroup', 'Author', 'A', TRUE);
end;
```

Nested Field Sections Field sections cannot be nested. For example, if you define a field section to start with `<TITLE>` and define another field section to start with `<FOO>`, the two sections *cannot* be nested as follows:

```
<TITLE> dog <FOO> cat </FOO> </TITLE>
```

To work with nested sections, define them as zone sections.

Repeated Field Sections Repeated field sections are allowed, but `WITHIN` queries treat them as a single section. The following is an example of repeated field section in a document:

```
<TITLE> cat </TITLE>
<TITLE> dog </TITLE>
```

The query *dog and cat within title* returns the document, even though these words occur in different sections.

To have `WITHIN` queries distinguish repeated sections, define them as zone sections.

Attribute Section

You can define attribute sections to query on XML attribute text. You can also have the system automatically define and index XML attributes for you.

See Also: ["XML Section Searching"](#) in this chapter.

Special Section

Special sections are not recognized by tags. Currently the only special sections supported are sentence and paragraph. This enables you to search for combination of words within sentences or paragraphs.

To add a special section, use the `CTX_DDL.ADD_SPECIAL_SECTION` procedure. For example, the following code enables searching within sentences within HTML documents:

```
begin
ctx_ddl.create_section_group('htmgroup', 'HTML_SECTION_GROUP');
ctx_ddl.add_special_section('htmgroup', 'SENTENCE');
end;
```

You can also add zone sections to the group to enable zone searching in addition to sentence searching. The following example adds the zone section `Headline` to the section group `htmgroup`:

```
begin
ctx_ddl.create_section_group('htmgroup', 'HTML_SECTION_GROUP');
ctx_ddl.add_special_section('htmgroup', 'SENTENCE');
ctx_ddl.add_zone_section('htmgroup', 'Headline', 'H1');
end;
```

HTML Section Searching

HTML has internal structure in the form of tagged text which you can use for section searching. For example, you can define a section called headings for the <H1> tag. This allows you to search for terms only within these tags across your document set.

To query, you use the WITHIN operator. Oracle returns all documents that contain your query term within the headings section. Thus, if you wanted to find all documents that contain the word oracle within headings, you issue the following query:

```
'oracle within headings'
```

Creating HTML Sections

The following code defines a section group called `htmgroup` of type `HTML_SECTION_GROUP`. It then creates a zone section in `htmgroup` called `headline` identified by the <H1> tag:

```
begin
ctx_ddl.create_section_group('htmgroup', 'HTML_SECTION_GROUP');
ctx_ddl.add_zone_section('htmgroup', 'heading', 'H1');
end;
```

You can then index your documents as follows:

```
create index myindex on docs(htmlfile) indextype is ctxsys.context
parameters('filter ctxsys.null_filter section group htmgroup');
```

After indexing with section group `htmgroup`, you can query within the heading section by issuing a query as follows:

```
'Oracle WITHIN heading'
```

Searching HTML Meta Tags

With HTML documents you can also create sections for NAME/CONTENT pairs in <META> tags. When you do so you can limit your searches to text within CONTENT.

Example: Creating Sections for <META>Tags

Consider an HTML document that has a META tag as follows:

```
<META NAME="author" CONTENT="ken">
```

To create a zone section that indexes all CONTENT attributes for the META tag whose NAME value is author:

```
begin
ctx_ddl.create_section_group('htmgroup', 'HTML_SECTION_GROUP');
ctx_ddl.add_zone_section('htmgroup', 'author', 'meta@author');
end
```

After indexing with section group htmgroup, you can query the document as follows:

```
'ken WITHIN author'
```

XML Section Searching

Like HTML documents, XML documents have tagged text which you can use to define blocks of text for section searching. The contents of a section can be searched on with the `WITHIN` or `INPATH` operators.

For XML searching, you can do the following:

- automatic sectioning
- attribute searching
- document type sensitive sections
- path section searching

Automatic Sectioning

You can set up your indexing operation to automatically create sections from XML documents using the section group `AUTO_SECTION_GROUP`. The system creates zone sections for XML tags. Attribute sections are created for the tags that have attributes and these sections named in the form `tag@attribute`.

For example, the following command creates the index `myindex` on a column containing the XML files using the `AUTO_SECTION_GROUP`:

```
CREATE INDEX myindex ON xmldocs(xmlfile) INDEXTYPE IS ctxsys.context PARAMETERS ('datastore ctxsys.default_datastore filter ctxsys.null_filter section group ctxsys.auto_section_group');
```

Attribute Searching

You can search XML attribute text in one of two ways:

- Create attribute sections with `CTX_DDL.ADD_ATTR_SECTION` and then index with the `XML_SECTION_GROUP`. If you use `AUTO_SECTION_GROUP` when you index, attribute sections are created automatically. You can query attribute sections with the `WITHIN` operator.
- Index with the `PATH_SECTION_GROUP` and query attribute text with the `INPATH` operator.

Creating Attribute Sections

Consider an XML file that defines the BOOK tag with a TITLE attribute as follows:

```
<BOOK TITLE="Tale of Two Cities">
  It was the best of times.
</BOOK>
```

To define the title attribute as an attribute section, create an XML_SECTION_GROUP and define the attribute section as follows:

```
begin
ctx_ddl.create_section_group('myxmlgroup', 'XML_SECTION_GROUP');
ctx_ddl.add_attr_section('myxmlgroup', 'booktitle', 'book@title');
end;
```

To index:

```
CREATE INDEX myindex ON xmldocs(xmlfile) INDEXTYPE IS ctxsys.context PARAMETERS
('datastore ctxsys.default_datastore filter ctxsys.null_filter section group
myxmlgroup');
```

You can query the XML attribute section *booktitle* as follows:

```
'Cities within booktitle'
```

Searching Attributes with the INPATH Operator

You can search attribute text with the INPATH operator. To do so, you must index your XML document set with the PATH_SECTION_GROUP.

See Also: ["Path Section Searching"](#) in this chapter.

Creating Document Type Sensitive Sections

You have an XML document set that contains the <book> tag declared for different document types. You want to create a distinct book section for each document type.

Assume that mydocname1 is declared as an XML document type (root element) as follows:

```
<!DOCTYPE mydocname1 ... [...
```

Within `mydocname1`, the element `<book>` is declared. For this tag, you can create a section named `mybooksec1` that is sensitive to the tag's document type as follows:

```
begin
  ctx_ddl.create_section_group('myxmlgroup', 'XML_SECTION_GROUP');
  ctx_ddl.add_zone_section('myxmlgroup', 'mybooksec1', 'mydocname1(book)');
end;
```

Assume that `mydocname2` is declared as another XML document type (root element) as follows:

```
<!DOCTYPE mydocname2 ... [...
```

Within `mydocname2`, the element `<book>` is declared. For this tag, you can create a section named `mybooksec2` that is sensitive to the tag's document type as follows:

```
begin
  ctx_ddl.create_section_group('myxmlgroup', 'XML_SECTION_GROUP');
  ctx_ddl.add_zone_section('myxmlgroup', 'mybooksec2', 'mydocname2(book)');
end;
```

To query within the section `mybooksec1`, use `WITHIN` as follows:

```
'oracle within mybooksec1'
```

Path Section Searching

XML documents can have parent-child tag structures such as the following:

```
<A> <B> <C> dog </C> </B </A>
```

In this example, tag `C` is a child of tag `B` which is a child of tag `A`.

With Oracle Text, you can do path searching with `PATH_SECTION_GROUP`. This section group allows you to specify direct parentage in queries, such as to find all documents that contain the term *dog* in element `C` which is a child of element `B` and so on.

With `PATH_SECTION_GROUP`, you can also perform attribute value searching and attribute equality testing.

The new operators associated with this feature are

- `INPATH`
- `HASPATH`

Creating Index with PATH_SECTION_GROUP

To enable path section searching, index your XML document set with PATH_SECTION_GROUP.

Create the preference:

```
begin
ctx_ddl.create_section_group('xmlpathgroup', 'PATH_SECTION_GROUP');
end;
```

Create the index:

```
CREATE INDEX myindex ON xmldocs(xmlfile) INDEXTYPE IS ctxsys.context PARAMETERS
('datastore ctxsys.default_datastore filter ctxsys.null_filter section group
xmlpathgroup');
```

When you create the index, you can use the INPATH and HASPATH operators.

Top-Level Tag Searching

To find all documents that contain the term *dog* in the top-level tag <A>:

```
dog INPATH (/A)
or
dog INPATH(A)
```

Any-Level Tag Searching

To find all documents that contain the term *dog* in the <A> tag at any level:

```
dog INPATH(//A)
```

This query finds the following documents:

```
<A>dog</A>
```

and

```
<A><B><C>dog</C></B></A>
```

Direct Parentage Searching

To find all documents that contain the term *dog* in a B element that is a direct child of a top-level A element:

```
dog INPATH(A/B)
```

This query finds the following XML document:

```
<A><B>My dog is friendly.</B></A>
```

but does not find:

```
<C><B>My dog is friendly.</B></C>
```

Tag Value Testing

You can test the value of tags. For example, the query:

```
dog INPATH(A[B="dog" ])
```

Finds the following document:

```
<A><B>dog</B></A>
```

But does not find:

```
<A><B>My dog is friendly.</B></A>
```

Attribute Searching

You can search the content of attributes. For example, the query:

```
dog INPATH(//A/@B)
```

Finds the document

```
<C><A B="snoop dog"> </A> </C>
```

Attribute Value Testing

You can test the value of attributes. For example, the query

```
California INPATH (//A[@B = "home address"])
```

Finds the document:

```
<A B="home address">San Francisco, California, USA</A>
```

But does not find:

```
<A B="work address">San Francisco, California, USA</A>
```

Path Testing

You can test if a path exists with the HASPATH operator. For example, the query:

```
HASPATH(A/B/C)
```

finds and returns a score of 100 for the document

```
<A><B><C>dog</C></B></A>
```

without the query having to reference *dog* at all.

Section Equality Testing with HASPATH

You can use the HASPATH operator to do section quality tests. For example, consider the following query:

```
dog INPATH A
```

finds

```
<A>dog</A>
```

but it also finds

```
<A>dog park</A>
```

To limit the query to the term *dog* and nothing else, you can use a section equality test with the HASPATH operator. For example,

```
HASPATH(A="dog" )
```

finds and returns a score of 100 only for the first document, and not the second.

See Also: *Oracle Text Reference* to learn more about using the INPATH and HASPATH operators.

Working With a Thesaurus

This chapter describes how to improve your query application with a thesaurus. The following topics are discussed in this chapter:

- [Overview of Thesauri](#)
- [Defining Thesaural Terms](#)
- [Using a Thesaurus in a Query Application](#)
- [About the Supplied Knowledge Base](#)

Overview of Thesauri

Users of your query application looking for information on a given topic might not know which words have been used in documents that refer to that topic.

Oracle Text enables you to create case-sensitive or case-insensitive thesauri which define synonym and hierarchical relationships between words and phrases. You can then retrieve documents that contain relevant text by expanding queries to include similar or related terms as defined in the thesaurus.

You can create a thesaurus and load it into the system.

Note: The Oracle Text thesauri formats and functionality are compliant with both the ISO-2788 and ANSI Z39.19 (1993) standards.

Thesaurus Creation and Maintenance

Thesauri and thesaurus entries can be created, modified, and deleted by all Oracle Text users with the CTXAPP role.

CTX_THES Package

To maintain and browse your thesaurus programatically, you can use the PL/SQL package, CTX_THES. With this package, you can browse terms and hierarchical relationships, add and delete terms, and add and remove thesaurus relations.

Thesaurus Operators

You can also use the thesaurus operators in the CONTAINS clause to expand query terms according to your loaded thesaurus. For example, you can use the SYN operator to expand a term such as *dog* to its synonyms as follows:

```
'syn(dog)'
```

ctxload Utility

The ctxload utility can be used for loading (creating) thesauri from a plain-text file into the thesaurus tables, as well as dumping thesauri from the tables into output (dump) files.

The thesaurus dump files created by ctxload can be printed out or used as input for other applications. The dump files can also be used to load a thesaurus into the

thesaurus tables. This can be useful for using an existing thesaurus as the basis for creating a new thesaurus.

Case-sensitive Thesauri

In a case-sensitive thesaurus, terms (words and phrases) are stored exactly as entered. For example, if a term is entered in mixed-case (using either the CTX_THES package or a thesaurus load file), the thesaurus stores the entry in mixed-case.

Note: To take full advantage of query expansions that result from a case-sensitive thesaurus, your index must also be case-sensitive.

When loading a thesaurus, you can specify that the thesaurus be loaded case-sensitive using the `-thescase` parameter.

When creating a thesaurus with `CTX_THES.CREATE_THESAURUS`, you can specify that the thesaurus created be case-sensitive.

In addition, when a case-sensitive thesaurus is specified in a query, the thesaurus lookup uses the query terms exactly as entered in the query. Therefore, queries that use case-sensitive thesauri allow for a higher level of precision in the query expansion, which helps lookup when and only when you have a case-sensitive index.

For example, a case-sensitive thesaurus is created with different entries for the distinct meanings of the terms *Turkey* (the country) and *turkey* (the type of bird). Using the thesaurus, a query for *Turkey* expands to include only the entries associated with *Turkey*.

Case-insensitive Thesauri

In a case-insensitive thesaurus, terms are stored in all-uppercase, regardless of the case in which they were entered.

The `ctxload` program loads a thesaurus case-insensitive by default.

When creating a thesaurus with `CTX_THES.CREATE_THESAURUS`, the thesaurus is created case-insensitive by default.

In addition, when a case-insensitive thesaurus is specified in a query, the query terms are converted to all-uppercase for thesaurus lookup. As a result, Oracle Text

is unable to distinguish between terms that have different meanings when they are in mixed-case.

For example, a case-insensitive thesaurus is created with different entries for the two distinct meanings of the term *TURKEY* (the country or the type of bird). Using the thesaurus, a query for either *Turkey* or *turkey* is converted to *TURKEY* for thesaurus lookup and then expanded to include all the entries associated with both meanings.

Default Thesaurus

If you do not specify a thesaurus by name in a query, by default, the thesaurus operators use a thesaurus named *DEFAULT*. However, Oracle Text does not provide a *DEFAULT* thesaurus.

As a result, if you want to use a default thesaurus for the thesaurus operators, you must create a thesaurus named *DEFAULT*. You can create the thesaurus through any of the thesaurus creation methods supported by Oracle Text:

- `CTX_THES.CREATE_THESAURUS` (PL/SQL)
- `ctxload`

See Also: *Oracle Text Reference* to learn more about using `ctxload` and the `CTX_THES` package.

Supplied Thesaurus

Although Oracle Text does not provide a default thesaurus, Oracle Text does supply a thesaurus, in the form of a `ctxload` load file, that can be used to create a general-purpose, English-language thesaurus.

The thesaurus load file can be used to create a default thesaurus for Oracle Text or it can be used as the basis for creating thesauri tailored to a specific subject or range of subjects.

See Also: *Oracle Text Reference* to learn more about using `ctxload` and the `CTX_THES` package.

Supplied Thesaurus Structure and Content

The supplied thesaurus is similar to a traditional thesaurus, such as Roget's Thesaurus, in that it provides a list of synonymous and semantically related terms.

The supplied thesaurus provides additional value by organizing the terms into a hierarchy that defines real-world, practical relationships between narrower terms and their broader terms.

Additionally, cross-references are established between terms in different areas of the hierarchy.

Supplied Thesaurus Location

The exact name and location of the thesaurus load file is operating system dependent; however, the file is generally named `dr0thsus` (with an appropriate extension for text files) and is generally located in the following directory structure:

```
<Oracle_home_directory>  
  <interMedia_Text_directory>  
    sample  
      thes
```

See Also: For more information about the directory structure for Oracle Text, see the Oracle9i installation documentation specific to your operating system.

Defining Thesaural Terms

You can create synonyms, related terms, and hierarchical relationships with a thesaurus. The following sections give examples.

Defining Synonyms

If you have a thesaurus of computer science terms, you might define a synonym for the term *XML* as *extensible markup language*. This allows queries on either of these terms to return the same documents.

```
XML
  SYN Extensible Markup Language
```

You can thus use the SYN operator to expand XML into its synonyms:

```
'SYN(XML)'
```

is expanded to:

```
'XML, Extensible Markup Language'
```

Defining Hierarchical Relations

If your document set is made up of news articles, you can use a thesaurus to define a hierarchy of geographical terms. Consider the following hierarchy that describes a geographical hierarchy for the U.S state of California:

```
California
  NT Northern California
    NT San Francisco
    NT San Jose
  NT Central Valley
    NT Fresno
  NT Southern California
    NT Los Angeles
```

You can thus use the NT operator to expand a query on California as follows:

```
'NT(California)'
```

expands to:

```
'California, Northern California, San Francisco, San Jose, Central Valley,
Fresno, Southern California, Los Angeles'
```

The resulting hitlist shows all documents related to the U.S. state of California regions and cities.

Using a Thesaurus in a Query Application

Defining a custom thesaurus allows you to process queries more intelligently. Since users of your application might not know which words represent a topic, you can define synonyms or narrower terms for likely query terms. You can use the thesaurus operators to expand your query into your thesaurus terms.

There are two ways to enhance your query application with a custom thesaurus so that you can process queries more intelligently:

- Load your custom thesaurus and issue queries with thesaurus operators
- Augment the knowledge base with your custom thesaurus (English only) and use the ABOUT operator to expand your query.

Each approach has its advantages and disadvantages.

Loading a Custom Thesaurus and Issuing Thesaural Queries

To build a custom thesaurus, follow these steps:

1. Create your thesaurus. See "[Defining Thesaural Terms](#)" in this chapter.
2. Load thesaurus with `ctxload`. For example, the following example imports a thesaurus named `tech_doc` from an import file named `tech_thesaurus.txt`:

```
ctxload -user jsmith/123abc -thes -name tech_doc -file tech_thesaurus.txt
```

3. Use THES operators to query. For example, you can find all documents that contain XML and its synonyms as defined in `tech_doc`:

```
'SYN(XML, tech_doc)'
```

Advantage

The advantage of using this method is that you can modify the thesaurus after indexing.

Limitations

This method requires you to use thesaurus expansion operators in your query. Long queries can cause extra overhead in the thesaurus expansion and slow your query down.

Augmenting Knowledge Base with Custom Thesaurus

You can add your custom thesaurus to a branch in the existing knowledge base. The knowledge base is a hierarchical tree of concepts used for theme indexing, ABOUT queries, and deriving themes for document services.

When you augment the existing knowledge base with your new thesaurus, you query with the ABOUT operator which implicitly expands to synonyms and narrower terms. You do not query with the thesaurus operators.

To augment the existing knowledge base with your custom thesaurus, follow these steps:

1. Create your custom thesaurus, linking new terms to existing knowledge base terms. See "[Defining Thesaural Terms](#)" and "[Linking New Terms to Existing Terms](#)".
2. Load thesaurus with `ctxload`. See "[Loading a Thesaurus with ctxload](#)".
3. Compile the loaded thesaurus with `ctxkbtcc` compiler. "[Compiling a Loaded Thesaurus](#)" later in this section.
4. Index your documents. By default the system creates a theme component to your index.
5. Use ABOUT operator to query. For example, to find all documents that are related to the term politics including any synonyms or narrower terms as defined in the knowledge base, issue the query:

```
'about(politics)'
```

Advantage

Compiling your custom thesaurus with the existing knowledge base before indexing allows for faster and simpler queries with the ABOUT operator. Document services can also take full advantage of the customized information for creating theme summaries and Gists.

Limitations

Use of the ABOUT operator requires a theme component in the index, which requires slightly more disk space. You must also define the thesaurus before indexing your documents. If you make any change to the thesuarus, you must recompile your thesaurus and re-index your documents.

Linking New Terms to Existing Terms

When adding terms to the knowledge base, Oracle recommends that new terms be linked to one of the categories in the knowledge base for best results in theme proving.

See Also: *Oracle Text Reference* for more information about the supplied English knowledge base.

If new terms are kept completely separate from existing categories, fewer themes from new terms will be proven. The result of this is poor precision and recall with ABOUT queries as well as poor quality of gists and theme highlighting.

You link new terms to existing terms by making an existing term the broader term for the new terms.

Example: Linking New Terms to Existing Terms You purchase a medical thesaurus `medthes` containing a hierarchy of medical terms. The four top terms in the thesaurus are the following:

- Anesthesia and Analgesia
- Anti-Allergic and Respiratory System Agents
- Anti-Inflammatory Agents, Antirheumatic Agents, and Inflammation Mediators
- Antineoplastic and Immunosuppressive Agents

To link these terms to the existing knowledge base, add the following entries to the medical thesaurus to map the new terms to the existing *health and medicine* branch:

```
health and medicine
  NT Anesthesia and Analgesia
  NT Anti-Allergic and Respiratory System Agents
  NT Anti-Inflammatory Agents, Antirheumatic Agents, and Inflammation Mediators
  NT Antineoplastic and Immunosuppressive Agents
```

Loading a Thesaurus with `ctxload`

Assuming the medical thesaurus is in a file called `med.thes`, you load the thesaurus as `medthes` with `ctxload` as follows:

```
ctxload -thes -thescase y -name medthes -file med.thes -user ctxsys/ctxsys
```


Compiling a Loaded Thesaurus

To link the loaded thesaurus `medthes` to the knowledge base, use `ctxkbtc` as follows:

```
ctxkbtc -user ctxsys/ctxsys -name medthes
```

About the Supplied Knowledge Base

Oracle Text supplies a knowledge base for English and French. The supplied knowledge contains the information used to perform theme analysis. Theme analysis includes theme indexing, ABOUT queries, and theme extraction with the CTX_DOC package.

The knowledge base is a hierarchical tree of concepts and categories. It has six main branches:

- science and technology
- business and economics
- government and military
- social environment
- geography
- abstract ideas and concepts

See Also: *Oracle Text Reference* for the breakdown of the category hierarchy.

The supplied knowledge base is like a thesaurus in that it is hierarchical and contains broader term, narrower term, and related term information. As such, you can improve the accuracy of theme analysis by augmenting the knowledge base with your industry-specific thesaurus by linking new terms to existing terms.

See Also: ["Augmenting Knowledge Base with Custom Thesaurus"](#) in this chapter.

You can also extend theme functionality to other languages by compiling a language-specific thesaurus into a knowledge base.

See Also: ["Adding a Language-Specific Knowledge Base"](#) in this chapter.

Knowledge Base Character Set

Knowledge bases can be in any single-byte character set. Supplied knowledge bases are in WE8ISO8859P1. You can store an extended knowledge base in another character set such as US7ASCII.

Adding a Language-Specific Knowledge Base

You can extend theme functionality to languages other than English or French by loading your own knowledge base for any single-byte whitespace delimited language, including Spanish.

Theme functionality includes theme indexing, ABOUT queries, theme highlighting, and the generation of themes, gists, and theme summaries with CTX_DOC.

You extend theme functionality by adding a user-defined knowledge base. For example, you can create a Spanish knowledge base from a Spanish thesaurus.

To load your language-specific knowledge base, follow these steps:

1. Load your custom thesaurus using `ctxload`.
2. Set `NLS_LANG` so that the language portion is the target language. The charset portion must be a single-byte character set.
3. Compile the loaded thesaurus using `ctxkbtbc`:

```
ctxkbtbc -user ctxsys/ctxsys -name my_lang_thes
```

This command compiles your language-specific knowledge base from the loaded thesaurus. To use this knowledge base for theme analysis during indexing and ABOUT queries, specify the `NLS_LANG` language as the `THEME_LANGUAGE` attribute value for the `BASIC_LEXER` preference.

Limitations

The following limitations hold for adding knowledge bases:

- Oracle supplies knowledge bases in English and French only. You must provide your own thesaurus for any other language.
- You can only add knowledge bases for languages with single-byte character sets. You cannot create a knowledge base for languages which can be expressed only in multi-byte character sets. If the database is a multi-byte universal character set, such as UTF-8, the `NLS_LANG` parameter must still be set to a compatible single-byte character set when compiling the thesaurus.
- Adding a knowledge base works best for whitespace delimited languages.
- You can have at most one knowledge base per NLS language.
- Obtaining hierarchical query feedback information such as broader terms, narrower terms and related terms does not work in languages other than English and French. In other languages, the knowledge bases are derived

entirely from your thesauri. In such cases, Oracle recommends that you obtain hierarchical information directly from your thesauri.

See Also: *Oracle Text Reference* for more information about theme indexing, ABOUT queries, using the CTX_DOC package, and the supplied English knowledge base.

Administration

This chapter describes Oracle Text administration. The following topics are covered:

- [Oracle Text Users and Roles](#)
- [DML Queue](#)
- [The CTX_OUTPUT Package](#)
- [Servers](#)
- [Administration Tool](#)

Oracle Text Users and Roles

While any user can create an Oracle Text index and issue a CONTAINS query, Oracle Text provides the CTXSYS user for administration and the CTXAPP role for application developers.

CTXSYS User

The CTXSYS user is created at install time. You administer Oracle Text users as this user.

CTXSYS can do the following:

- Modify system-defined preferences
- Drop and modify other user preferences
- Call procedures in the CTX_ADM PL/SQL package to set system-parameters
- Query all system-defined views
- Perform all the tasks of a user with the CTXAPP role

CTXAPP Role

The CTXAPP role is a system-defined role that enables users to do the following:

- Create and delete Oracle Text preferences
- Use the Oracle Text PL/SQL packages

Any user can create an Oracle Text index and issue a Text query. The CTXAPP role allows users create preferences and use the PL/SQL packages.

Granting Roles and Privileges to Users

The system uses the standard SQL model for granting roles to users. To grant a Text role to a user, use the GRANT statement.

In addition, to allow application developers to call procedures in the Oracle Text PL/SQL packages, you must explicitly grant to each user EXECUTE privileges for the Oracle Text package.

DML Queue

When there are inserts, updates, or deletes to documents in your base table, the DML queue stores the requests for documents waiting to be indexed. When you synchronize the index with `CTX_DDL.SYNC_INDEX`, requests are removed from this queue.

Pending DML requests can be queried with the `CTX_PENDING` and `CTX_USER_PENDING` views.

DML errors can be queried with the `CTX_INDEX_ERRORS` or `CTX_USER_INDEX_ERRORS` view.

See Also: *Oracle Text Reference* for more information about these views.

The CTX_OUTPUT Package

Use the CTX_OUTPUT PL/SQL package to log indexing and document service requests.

See Also: *Oracle Text Reference* for more information about this package.

Servers

You index documents and issue queries with standard SQL. No server is needed for performing batch DML. You can synchronize the CONTEXT index with the CTX_DDL.SYNC_INDEX procedure.

See Also: For more information about indexing and index synchronization, see [Chapter 2, "Indexing"](#).

Administration Tool

The Oracle Text Manager is a Java application integrated with the Oracle Enterprise Manager, which is available on a separate CD.

The Text Manager enables administrators to create preferences, stoplists, sections, and indexes. This tool also enables administrators to perform DML.

See Also: for more information about the Oracle Text Manager, see the online help shipped with this tool.

CONTEXT Query Application

This appendix describes how to build a simple web-search application using the CONTEXT index type. The following topic is covered:

- [Web Query Application Overview](#)
- [The PSP Web Application](#)
- [Web Application Sample Code](#)

Web Query Application Overview

A common use of Oracle Text is to index HTML files on web sites and provide search capabilities to users. The sample application in this Appendix indexes a set of HTML files stored in the database and uses a web server connected to Oracle to provide the search service.

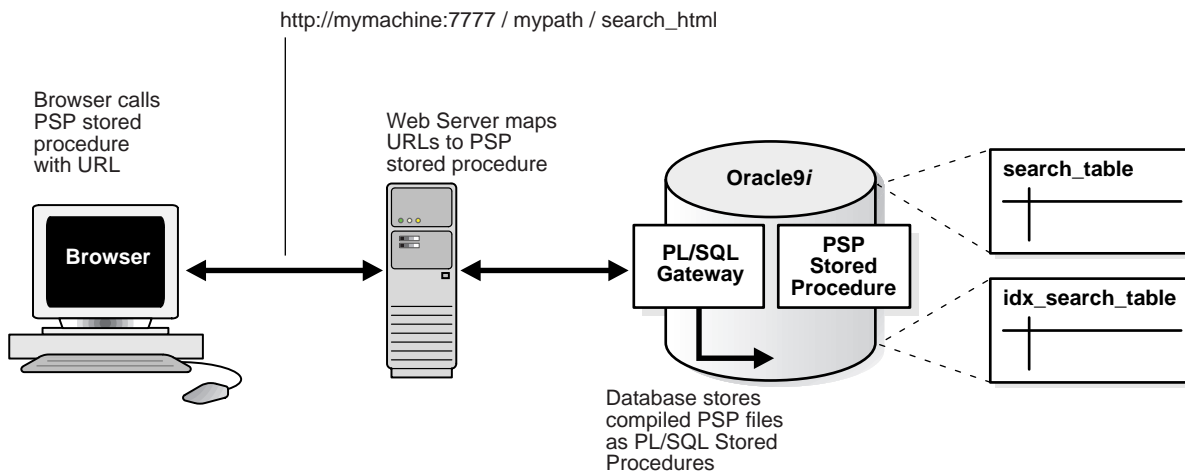
There are two versions of this application. One that uses PL/SQL Server Pages (PSP) and one that uses Java Server Pages (JSP). This appendix describes the PSP application. You can view and download both the PSP and JSP application code at the Oracle Technology Network web site:

<http://technet.oracle.com/products/text>

The PSP Web Application

This application is based on PL/SQL server pages. [Figure A-1](#) illustrates how the browser calls the PSP stored procedure on Oracle9i via a web server.

Figure A-1



Web Application Prerequisites

This application has the following requirements:

- Your Oracle database (version 8.1.6 or higher) is up and running.
- You have the Oracle PL/SQL gateway running
- You have a web server such as Apache up and running and correctly configured to send requests to the Oracle9i server.

Building the Web Application

This section describes how to build the web application.

Step 1 Create your Text Table

You must create a text table to store your html files. This example creates a table called `search_table` as follows:

```
create table search_table (tk numeric primary key, title varchar2(2000), text
clob);
```

Step 2 Load HTML Documents into Table Using SQL*Loader

You must load the text table with the HTML files. This example uses the control file [loader.ctl](#) to load the files named in [loader.dat](#). The SQL*Loader command is as follows:

```
% sqlldr userid=scott/tiger control=loader.ctl
```

Step 3 Create the CONTEXT index

Index the HTML files by creating a CONTEXT index on the text column as follows. Since we are indexing HTML, this example uses the NULL_FILTER preference type for no filtering and uses the HTML_SECTION_GROUP type:

```
create index idx_search_table on search_table(text)
  indextype is ctxsys.context parameters
  ('filter ctxsys.null_filter section group CTXSYS.HTML_SECTION_GROUP');
```

Step 4 Compile search_htmlservices Package in Oracle9i

The application must present selected documents to the user. To do so, Oracle must read the documents from the CLOB in `search_table` and output the result for viewing. This is done by calling procedures in the `search_htmlservices` package. The file [search_htmlservices.sql](#) must be compiled. You can do this at the SQL*Plus prompt:

```
SQL> @search_htmlservices.sql
```

```
Package created.
```

Step 5 Compile the search_html PSP page with loadpsp

The search page is invoked by calling [search_html.psp](#) from a browser. You compile search_html in Oracle9i with the loadpsp command-line program:

```
% loadpsp -replace -user scott/tiger search_html.psp
"search_html.psp": procedure "search_html" created.
```

See Also: *Oracle9i Application Developer's Guide - Fundamentals* for more information about using PSP.

Step 6 Configure Your Web Server

You must configure your web server to accept client PSP requests as a URL. Your web server forwards these requests to the Oracle9i server and returns server output to the browser. Refer to [Figure A-1](#).

You can use the Oracle WebDB 2.x web listener or Oracle iAS which includes the Apache web server. See your web server documentation for more information.

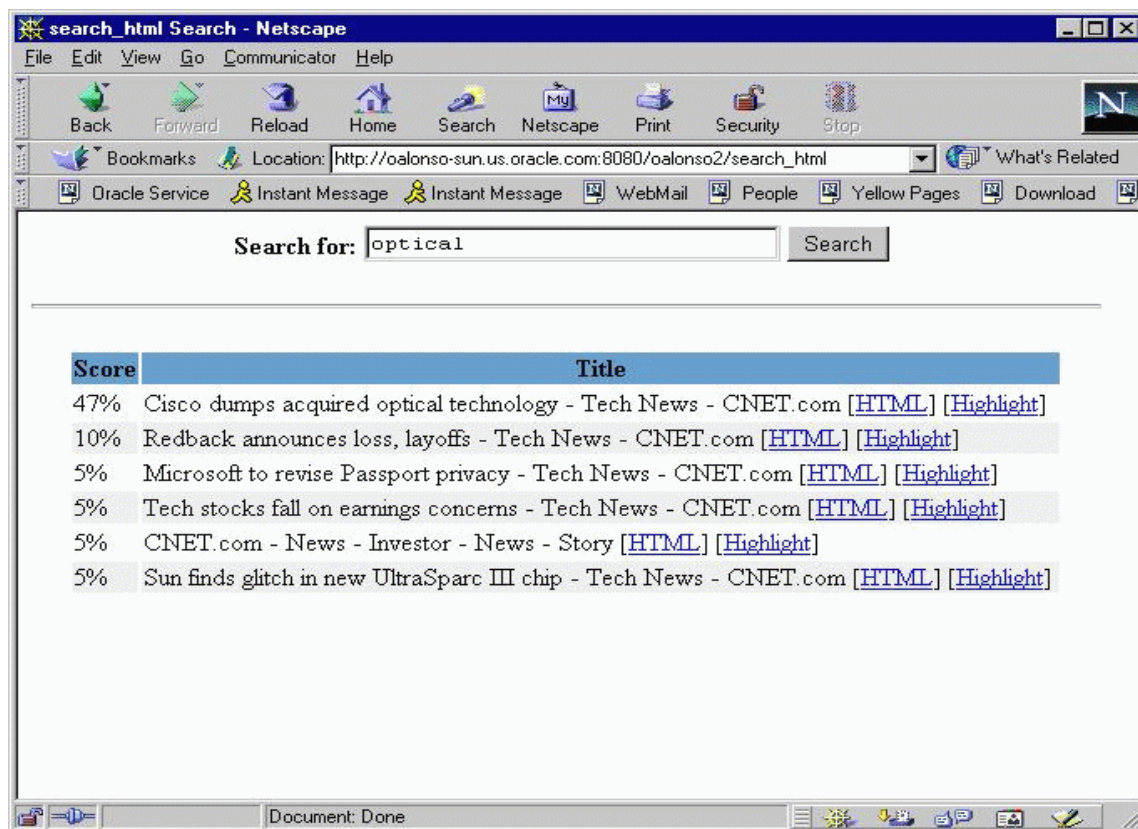
Step 7 Issue Query from Browser

You can access the query application from a browser using a URL. You configure the URL with your web server. An example URL might look like:

```
http://mymachine:7777/mypath/search_html
```

The application displays a query entry box in your browser and returns the query results as a list of HTML links. See [Figure A-2, "Screen shot of Web Query Application"](#).

Figure A-2 Screen shot of Web Query Application



Web Application Sample Code

This section lists the code used to build the example web application. It includes the following files:

- [loader.ctl](#)
- [loader.dat](#)
- [search_htmlservices.sql](#)
- [search_html.psp](#)

See Also: <http://technet.oracle.com/products/text/>

loader.ctl

```
LOAD DATA
  INFILE 'loader.dat'
  INTO TABLE search_table
  REPLACE
  FIELDS TERMINATED BY ';'
  (tk          INTEGER,
   title       CHAR,
   text_file   FILLER CHAR,
   text        LOBFILE(text_file) TERMINATED BY EOF)
```


loader.dat

```
1; Sun finds glitch in new UltraSparc III chip;0-1003-200-5507959.html
2; Redback announces loss, layoffs ;0-1004-200-5424681.html
3; Cisco dumps acquired optical technology ;0-1004-200-5510096.html
4; Microsoft to revise Passport privacy ;0-1005-200-5508903.html
5; Tech stocks fall on earnings concerns;0-1007-200-5506210.html
6; CNET.com - News - Investor - News - Story ;0-9900-1028-5510548-0.html
7; Chicago Tribune JUSTICES HEAR ARGUMENTS ;0_2669_SAV-0103290318_FF.html
8; Massive new effort to combat African AIDS is planned ;WEST04.html
9; U.S. Had Biggest Growth in 1990s ;census_2000.html
10; Congress Discusses Napster Issues ;congress_napster.html
11; Washington And China Face Off in Spy Plane Drama ;crash_china_dc_35.html
12; American Arrive To Study in Cuba ;cuba_us_medical_students_1.html
13; Hubble Spots Most-Distant Supernova ;distant_supernova.html
14; Survey: U.S. Has 90 Percent Chance of Recession;economy_forecast_dc_1.html
15; House Votes To Repeal Estate Tax ;estate_tax.html
16; EU Condemns Bush on Global Warming ;eu_global_warming.html
17; Foot-and-Mouth Vaccinations on Hold ;foot_and_mouth.html
18; Foot-and-Mouth Vaccinations on Hold ;foot_and_mouth_7.html
19; Cancer Research Project Links Millions of PCs ;health_cancer_dc_1.html
20; Company Says Early HIV Vaccine Data Are Promising ;hiv.html
21; Yahoo! Sports: SOW - Maradona Faces New Paternity Suit ;maradona.html
22; Israel, Palestinians Hold High-Level Talks ;mideast_leadall_dc.html
23; Evidence Mounts Against Milosevic ;milosevic_slain_rivals.html
24; Philippines Files Charges Against Estrada ;philippines_estrada_dc.html
25; Power Woes Affecting Calif. Economy ;power_woes.html
26; Dissidents Ask UN Rights Body to Condemn China ;rights_china_dc_2.html
27; South Africa to Act on Basis HIV Causes AIDS ;safrica_aids_dc_1.html
28; Shaggy Found Inspiration For Success In Jamaica ;shaggy_found.html
29; Solar Flare Eruptions Likely ;solar_flare.html
30; Plane Crash Kills Sudanese Officers ;sudan_plane_crash.html
31; SOUNDSCAN REPORT: Recipe for An Aspiring Top Ten;urban_groove_1.html
```

search_htmlservices.sql

```
set define off

create or replace package search_htmlServices as

    procedure showHTMLDoc (p_id in numeric);

    procedure showDoc (p_id in numeric, p_query in varchar2);

end;
/
show errors;

create or replace package body search_htmlServices as

    procedure showHTMLDoc (p_id in numeric) is
        v_clob_selected  CLOB;
        v_read_amount    integer;
        v_read_offset    integer;
        v_buffer          varchar2(32767);
    begin

        select text into v_clob_selected from search_table where tk = p_id;
        v_read_amount := 32767;
        v_read_offset := 1;
    begin
    loop
        dbms_lob.read(v_clob_selected,v_read_amount,v_read_offset,v_buffer);
        http.print(v_buffer);
        v_read_offset := v_read_offset + v_read_amount;
        v_read_amount := 32767;
    end loop;
    exception
    when no_data_found then
        null;
    end;
end showHTMLDoc;

    procedure showDoc (p_id in numeric, p_query in varchar2) is

        v_clob_selected  CLOB;
```

```
v_read_amount    integer;
v_read_offset    integer;
v_buffer         varchar2(32767);
v_query         varchar(2000);
v_cursor        integer;

begin
  http.p('<html><title>HTML version with highlighted terms</title>');
  http.p('<body bgcolor="#ffffff">');
  http.p('<b>HTML version with highlighted terms</b>');

  begin
    ctx_doc.markup (index_name => 'search_table_gen_index',
                   textkey    => p_id,
                   text_query => p_query,
                   restab     => v_clob_selected,
                   starttag   => '<i><font color=red>',
                   endtag     => '</font></i>');

    v_read_amount := 32767;
    v_read_offset := 1;
    begin
      loop
        dbms_lob.read(v_clob_selected,v_read_amount,v_read_offset,v_buffer);
        http.print(v_buffer);
        v_read_offset := v_read_offset + v_read_amount;
        v_read_amount := 32767;
      end loop;
    exception
      when no_data_found then
        null;
      end;

    exception
      when others then
        null; --showHTMLdoc(p_id);
      end;
  end showDoc;
end;
/
show errors

set define on
```

search_html.psp

```
<%@ plsql procedure="search_html" %>
<%@ plsql parameter="query" default="null" %>
<%! v_results numeric := 0; %>

<html>
<head>
  <title>search_html Search </title>
</head>
<body>

<%

If query is null Then

  -- This part of the script allows a person
  -- to enter data on an HTML form.
  %>

  <center>
    <form method=post action="search_html">
      <b>Search for: </b>
      <input type=text name="query" size=30>&nbsp;
      <input type=submit value=Search>
    </center>
  <hr>

  <%
  Else
  %>

  <p>
  <%!
    color varchar2(6) := 'ffffff';
  %>

  <center>
    <form method=post action="search_html">
      <b>Search for: </b>
      <input type=text name="query" size=30 value="<%= query %>">
      <input type=submit value=Search>
    </form>
  </center>
  <hr>
```

```

<p>

<%
  -- select statement
  for doc in (
    select /*+ FIRST_ROWS */ rowid, tk, title, score(1) scr
    from search_table
    where contains(text, query,1) >0
    order by score(1) desc
  )
  loop
    v_results := v_results + 1;
    if v_results = 1 then

%>

        <center>
          <table border="0">
            <tr bgcolor="#6699CC">
              <th>Score</th>
              <th>Title</th>
            </tr>

<%
  end if; %>
      <tr bgcolor="#<%= color %>">
        <td> <%= doc.scr %>% </td>
        <td> <%= doc.title %>
          [<a href="search_htmlServices.showHTMLDoc?p_id=<%= doc.tk
%>">HTML</a>]
          [<a href="search_htmlServices.showDoc?p_id=<%= doc.tk %>&p_query=<%=
query %>">Highlight</a>]
        </td>
      </tr>

<%
  if (color = 'ffffff') then
    color := 'eeeeee';
  else
    color := 'ffffff';
  end if;

  end loop;
%>

</table>

```

```
        </center>
    <%
        end if;
    %>
</body></html>
```

Index

A

ABOUT query, 3-13
 adding for your language, 7-13
 case-sensitivity, 3-11
 definition, 3-8
 example, 1-17
accents
 indexing characters with, 2-15
ACCUM operator, 3-14
ADD_STOPCLASS procedure, 2-26
ADD_STOPTHEME procedure, 2-26
ADD_STOPWORD procedure, 2-25, 2-26
ADD_SUB_LEXER procedure
 example, 2-23
administration tool, 8-6
ALL_ROWS hint
 better response time, 5-6
ALTER INDEX command
 rebuilding index, 2-35
 resuming failed index, 2-34
alternate spelling, 2-16
AND operator, 3-14
application
 sample, A-1
attribute
 searching XML, 6-12
attribute sections, 6-8
AUTO_SECTION_GROUP object, 6-3
automatic sections, 6-12

B

background DML, 8-5

base-letter conversion, 2-15
BASIC_LEXER, 2-13
BASIC_SECTION_GROUP object, 6-2
BFILE column, 1-7
 indexing, 1-10, 2-27
BINARY
 format column value, 2-12
BLOB column, 1-7
 indexing, 1-10, 2-27
blocking operations
 tuning queries with, 5-9
bypassing rows, 2-12

C

case-sensitive
 ABOUT query, 3-11
 indexing, 2-15
 queries, 3-10
 thesaurus, 7-3
CATSEARCH, 3-4
 creating index for, 2-30
 operators, 3-18
 SQL example, 3-4
 structured query, 3-5
CHAR column, 1-7
character set
 indexing, 2-12
 indexing mixed, 2-13
character set column, 1-7
charset column, 2-13
CHARSET_FILTER, 2-5, 2-13
Chinese indexing, 2-17
CHINESE_VGRAM_LEXER, 2-17

CHOOSE hint
 better response time, 5-6
CLOB column, 1-7
 indexing, 1-10, 2-27
column types
 supported for indexing, 1-7
composite words
 indexing, 2-16
concept query, See **ABOUT**
CONTAINS
 operators, 3-13
 PL/SQL example, 3-3
 query, 3-2
 SQL example, 3-2
 structured query, 3-3
CONTEXT index, 1-2
 about, 1-10, 2-9
 creating, 1-10, 2-20, 2-27
 customizing, 1-11
 HTML example, 2-28, A-3
counting hits, 3-21
CREATE INDEX command, 2-27
CREATE_STOPLIST procedure, 2-25, 2-26
CTX_DDL.SYNC_INDEX procedure, 2-37
CTX_DOC package, 4-2
CTX_INDEX_ERRORS view, 2-34, 8-3
CTX_PENDING view, 8-3
CTX_THES package
 about, 7-2
CTX_USER_INDEX_ERRORS view, 2-34, 8-3
CTX_USER_PENDING view, 8-3
CTXAPP role, 8-2
CTXCAT index, 1-2, 1-12
 about, 1-10, 2-9
 example, 2-29
ctxkbt
 example, 7-11
ctxload
 load thesaurus example, 7-2, 7-8, 7-10
CTXRULE index, 1-2, 1-12
 about, 1-10, 2-10
 creating, 2-32
CTXSYS user, 8-2

D

data storage
 index default, 1-10, 2-27
 preference example, 2-22
datastore
 about, 2-4, 2-20
DATE column, 1-10, 2-27
DBMS_JOB.SUBMIT procedure, 2-37
default thesaurus, 7-4
defaults
 index, 1-10, 2-27
DETAIL_DATASTORE, 1-6
 about, 2-11
diacritical marks
 characters with, 2-15
DIRECT_DATASTORE, 1-6
 about, 2-11
 example, 2-21
DML
 view pending, 2-36
DML processing, 1-12
 background, 8-5
DML queue, 8-3
document classification, 2-32
 about, 1-2
document formats
 filtering, 2-11
 supported, 1-3, 1-7
document hit count
 presenting, 1-22
document invalidation, 2-38
document loading
 methods, 1-8
document presentation
 about, 1-23
document sections, 2-25
document services
 about, 1-23
DOMAIN_INDEX_NO_SORT hint
 better throughput example, 5-8
DOMAIN_INDEX_SORT hint
 better response time example, 5-6
drjobdml.sql script, 2-37
DROP INDEX command, 2-34

DROP_STOPLIST procedure, 2-26
dropping an index, 2-34

E

EQUIV operator, 3-14

errors

 DML, 8-3

 viewing, 2-34

explain plan, 3-11

extensible query optimizer, 5-2

F

feedback

 query, 3-11

field section

 definition, 6-7

 nested, 6-8

 repeated, 6-8

 visible and invisible, 6-7

file paths

 storing, 1-6

FILE_DATASTORE, 2-4

 about, 1-6, 2-11

 example, 2-22

filter

 about, 2-5, 2-20

FILTER procedure, 4-3

filtering

 custom, 2-12

 index default, 1-11, 2-27

 to plain text and HTML, 1-23

filtering documents, 2-11

 to HTML and plain text, 4-3

FIRST_ROWS hint, 3-19

 better response time example, 5-5

 better throughput example, 5-8

 example, 1-17

format column, 1-7, 2-12

formats

 filtering, 2-11

 supported, 1-7

fragmentation of index, 2-37

full themes

 obtaining, 4-5

fuzzy matching, 2-17

 default, 1-11, 2-28

fuzzy operator, 3-15

G

garbage collection, 2-38

gist

 definition, 4-4

 example, 4-5

GIST procedure, 4-5

granting roles, 8-2

H

HASPATH operator, 6-14

 examples, 6-17

HFEEDBACK procedure, 3-11

HIGHLIGHT procedure, 4-2

highlighting

 about, 1-23

 overview, 4-2

highlighting text, 4-2

highlighting themes, 4-2

hit count, 3-21

hitlist

 presenting, 1-20

HTML

 filtering to, 1-23, 4-3

 indexing, 2-22, 6-2

 indexing example, A-3

 searching META tags, 6-10

 zone section example, 2-25, 6-10

HTML_SECTION_GROUP object, 2-25, 6-2, 6-10

 with NULL_FILTER, 2-22, A-3

I

IGNORE

 format column value, 2-12

index

 about, 2-2

 creating, 2-20, 2-27

 dropping, 2-34

- multiple, 2-7
- optimizing, 2-37, 2-38
- rebuilding, 2-35
- structure, 2-2, 2-37
- synchronizing, 2-36, 8-5
- index defaults
 - general, 1-10, 2-27
- index engine
 - about, 2-5
- index errors
 - viewing, 2-34
- index fragmentation, 2-37
- index maintenance, 1-12, 2-34
- index types
 - about, 1-9
 - choosing, 2-9
- indexing
 - about, 1-9
 - bypassing rows, 2-12
 - considerations, 2-8
 - limitations, 2-7
 - overview of process, 2-4
 - parallel, 2-6
 - resuming failed, 2-34
 - special characters, 2-13
- indexing views, 2-7
- INPATH operator, 6-14
 - examples, 6-15
- INSERT statement
 - load text example, 1-8
- INSO_FILTER, 2-5, 2-11, 2-13

J

- Japanese indexing, 2-17
- JAPANESE_LEXER, 2-17

K

- knowledge base
 - about, 7-12
 - augmenting, 7-9
 - linking new terms, 7-10
 - supported character set, 7-12
 - user-defined, 7-13

- Korean indexing, 2-17
- KOREAN_MORP_LEXER, 2-17

L

- language
 - default setting for indexing, 1-11, 2-27
- language specific features, 2-15
- languages
 - indexing, 2-13
- language-specific knowledge base, 7-13
- lexer
 - about, 2-5, 2-20
- list of themes
 - definition, 4-4
 - obtaining, 4-4
- loading text
 - about, 1-5
 - SQL INSERT example, 1-8
- LOB columns
 - indexing, 1-10, 2-27
- location of text, 2-10
- logical operators, 3-13

M

- maintaining the index, 2-34
- marked-up document
 - obtaining, 4-3
- MARKUP procedure, 4-3
- MATCHES
 - about, 3-5
 - PL/SQL example, 2-33, 3-6
 - SQL example, 3-5
- META tag
 - creating zone section for, 6-11
- mixed formats
 - filtering, 2-12
- MULTI_COLUMN_DATASTORE, 1-6
 - about, 2-11
 - example, 2-21
- MULTI_LEXER, 2-13
 - example, 2-23
- multi-language columns
 - indexing, 2-13

multi-language stoplist
 about, 2-26
multiple indexes, 2-7

N

NCLOB column, 1-10, 2-27
NEAR operator, 3-15
nested zone sections, 6-6
NESTED_DATASTORE, 1-6
 about, 2-11
NEWS_SECTION_GROUP object, 6-3
NOT operator, 3-14
NULL_FILTER, 2-5
 example, 2-22, A-3
NULL_SECTION_GROUP object, 6-2
NUMBER column, 1-10, 2-27

O

offset information
 highlight, 4-2
operators
 CATSEARCH, 3-18
 CONTAINS, 3-13
 logical, 3-13
 thesaurus, 7-2
optimizing index, 2-37
 example, 2-38
 single token, 2-38
optimizing queries, 3-19, 5-2
 response time, 1-17, 5-5
 statistics, 5-2
 throughput, 5-8
 with blocking operations, 5-9
OR operator, 3-14
Oracle Enterprise Manager, 8-6
Oracle9i Text Manager, 8-6

P

parallel indexing, 2-6
paramstring for CREATE INDEX, 2-27
path section searching, 6-14
PATH_SECTION_GROUP

 example, 6-15
pending DML
 viewing, 2-36
pending updates, 8-3
phrase query, 3-7
plain text
 filtering to, 4-3
 indexing with NULL_FILTER, 2-22
plain text filtering, 1-23
PL/SQL functions
 calling in contains, 3-16
preferences
 creating (examples), 2-21
 creating with admin tool, 8-6
 dropping, 2-35
presenting hitlist, 1-20
printjoins character, 2-14
PROCEDURE_FILTER, 2-12
PSP application, A-2

Q

query
 ABOUT, 3-13
 about, 1-16
 blocking operations, 5-9
 case-sensitive, 3-10
 CATSEARCH, 3-4
 CONTAINS, 3-2
 counting hits, 3-21
 MATCHES, 3-5
 optimizing for throughput, 5-8
 overview, 3-2
query application
 prerequisites, 1-4
 sample, 1-14
query example, 1-16
query explain plan, 3-11
query expressions, 3-9
query features, 1-18
query feedback, 3-11
query optimization, 3-19
query tuning, 5-5
queue
 DML, 8-3

R

- rebuilding an index, 2-35
- REMOVE_SQE procedure, 3-16
- REMOVE_STOPCLASS procedure, 2-26
- REMOVE_STOPTHEME procedure, 2-26
- REMOVE_STOPWORD procedure, 2-25, 2-26
- response time
 - improving, 5-5
 - optimizing for, 1-17, 3-19
- result buffer size
 - increasing, 5-9
- resuming failed index, 2-34
- roles
 - granting, 8-2
 - system-defined, 8-2

S

- score
 - presenting, 1-22
- section
 - attribute, 6-8
 - field, 6-7
 - HTML example, 2-25
 - nested, 6-6
 - overlapping, 6-6
 - repeated zone, 6-6
 - special, 6-9
 - zone, 6-5
- section group
 - about, 2-20
 - creating with admin tool, 8-6
- section searching
 - about, 1-18, 6-2
 - enabling, 6-2
 - HTML, 6-10
- sectioner
 - about, 2-5
- sectioning
 - automatic, 6-12
 - path, 6-14
- single themes
 - obtaining, 4-5
- skipjoins character, 2-14

- SORT_AREA_SIZE parameter, 5-9
- special characters
 - indexing, 2-13
- special sections, 6-9
- spelling
 - alternate, 2-16
- SQE operator, 3-15
- statistics
 - optimizing with, 5-2
- stem operator, 2-17, 3-15
- stemming
 - default, 1-11, 2-28
- stopclass, 2-26
- stoplist, 2-25
 - about, 2-20
 - creating with admin tool, 8-6
 - default, 1-11, 2-28
 - multi-language, 2-18, 2-26
 - PL/SQL procedures, 2-26
- stoptheme, 2-26
 - about, 2-18
 - definition, 3-9
- stopword, 2-25, 2-26
 - about, 2-18, 3-8
 - case-sensitive, 3-10
- storage
 - about, 2-20
- STORE_SQE procedure, 3-15, 3-16
- stored query expressions, 3-15
- storing text, 2-10
 - about, 1-6
- structure of index, 2-37
- structured field searching
 - about, 1-17
- structured fields
 - presenting in application, 1-22
- structured query
 - example, 2-29
- SYN operator, 7-6
- SYNC_INDEX procedure, 2-37
- synchronizing index, 1-12, 2-36, 8-5
- synonyms
 - defining, 7-6

T

TEXT

- format column value, 2-12
- text column
 - supported types, 1-7
- text highlighting, 4-2
- Text Manager tool, 8-6
- text query applications
 - about, 1-2
- text storage, 2-10
- theme capabilities
 - overview, 1-3
- theme functionality
 - adding, 7-13
- theme highlighting, 4-2
- theme query, See ABOUT
- theme summary
 - definition, 4-4
- themes
 - indexing, 2-15
- THEMES procedure, 4-4
- thesaural queries
 - about, 1-18
- thesaurus
 - about, 7-2
 - adding to knowledge base, 7-9
 - case-sensitive, 7-3
 - DEFAULT, 7-4
 - default, 7-4
 - defining terms, 7-6
 - hierarchical relations, 7-6
 - loading custom, 7-8
 - operators, 7-2
 - supplied, 7-4
 - using in application, 7-8
- thesaurus operator, 3-15
- throughput
 - improving query, 5-8
- tildes
 - indexing characters with, 2-15
- tuning queries
 - for response time, 1-17, 5-5
 - for throughput, 5-8
 - increasing result buffer size, 5-9

- with statistics, 5-2

U

- umlauts
 - indexing characters with, 2-15
- URL_DATASTORE
 - about, 2-11
 - example, 2-22
- URLs
 - storing, 1-6
- user
 - system-defined, 8-2
- USER_DATASTORE, 2-7
 - about, 2-11
- USER_FILTER, 2-12

V

- VARCHAR2 column, 1-7
- views
 - indexing, 2-7

W

- wildcard operator, 3-15
- WITHIN operator, 2-25
- word query, 3-7
 - case-sensitivity, 3-10
 - example, 1-16
- wordlist
 - about, 2-20

X

- XML documents
 - attribute searching, 6-12
 - doctype sensitive sections, 6-13
 - indexing, 6-3
 - section searching, 6-12
- XML_SECTION_GROUP object, 6-2

Z

- zone section
 - definition, 6-5

nested, 6-6
overlapping, 6-6
repeating, 6-6