

Oracle® Ultra Search

User's Guide

10g (9.0.4)

Part No. B10896-01

September 2003

Oracle Ultra Search User's Guide 10g (9.0.4)

Part No. B10896-01

Copyright © 2002, 2003 Oracle Corporation. All rights reserved.

Primary Author: Michele Cyran

Contributors: Sandeepan Banerjee, Stefan Buchta, Eddy Chee, Chung-Ho Chen, Will Chin, Jack Chung, Ray Hachem, Cindy Hsin, Hassan Karraby, Yasuhiro Matsuda, Colin McGregor, Valarie Moore, Visar Nimani, Steve Yang, David Zhang

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, except to the extent required to obtain interoperability with other independently created software or as specified by law, 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 Oracle9i, Oracle8i, PL/SQL, Oracle Store, SQL*Plus are trademarks or registered trademarks of Oracle Corporation. Other names may be trademarks of their respective owners.

Contents

Send Us Your Comments	xiii
Preface.....	xv
Audience	xv
Organization	xv
Related Documentation	xvii
Conventions.....	xviii
Documentation Accessibility	xxii
What's New in Ultra Search?.....	xxv
Ultra Search Release Information.....	xxviii
1 Introduction to Ultra Search	
Overview of Ultra Search.....	1-2
Ultra Search Components	1-2
Ultra Search Crawler.....	1-2
Ultra Search Backend	1-3
Ultra Search Administration Tool.....	1-3
Ultra Search APIs and Sample Applications	1-3
Ultra Search Features	1-4
Integration with Oracle Application Server	1-5
Extensible Crawler and Crawler Agents.....	1-5
Federated Search	1-6
Secure Search	1-6

Dependency on Oracle XML DB	1-7
Sample Query Applications	1-7
Sample Search Portlet.....	1-8
Query API.....	1-8
URL Rewrite	1-9
Robots Exclusions	1-9
Display URL Support.....	1-9
Document and Search Attributes	1-9
Metadata Loader.....	1-10
Document Relevancy Boosting.....	1-10
Data Harvesting Mode.....	1-11
Instance Snapshot Support.....	1-11
Integration with Oracle Internet Directory	1-11
Ultra Search Administration Groups in OID.....	1-11
Authorization of the Administration Privileges	1-12
Single Sign-On Authentication	1-12
Query Syntax Expansion	1-13
Ultra Search System Configuration.....	1-13

2 Installing and Configuring Ultra Search

Ultra Search Requirements	2-2
Ultra Search Conventions.....	2-2
Ultra Search System Requirements.....	2-2
Hardware Requirements	2-3
Software Requirements.....	2-4
Installing the Ultra Search Backend.....	2-4
Installing and Configuring the Ultra Search Backend	2-4
Configure a Secure Ultra Search Installation.....	2-6
Installing the Backend on an Existing Database or Metadata Repository	2-9
Database Requirements	2-9
Installing the Ultra Search Backend on an Existing Oracle9i Database with OPCA	2-9
Configuring the Default Ultra Search Instance	2-11
Installing the Ultra Search Middle Tier on Web Server Hosts	2-11
Web Applications Concepts.....	2-12
Browser Requirements.....	2-13

Installing the Middle Tier with the Oracle Database Release.....	2-13
Installing the Middle Tier with the Oracle Application Server Release.....	2-14
Configuring the Middle Tier with Oracle HTTP Server and OC4J	2-15
Configuring the Administration Tool with Single Sign-On Server	2-18
Deploying the Ultra Search EAR File on a Third Party Middle Tier	2-19
Editing the data-sources.xml File.....	2-22
Editing the ultrasearch.properties File.....	2-24
Starting the Web Server.....	2-25
Testing the Ultra Search Administration Tool	2-25
Testing the Ultra Search Sample Query Applications	2-26
Installing the Backend on Remote Crawler Hosts	2-27
Installing the Backend on Remote Crawler Hosts.....	2-27
Configuring the Backend on Remote Crawler Hosts.....	2-28
Unregistering a Remote Crawler.....	2-30
Configuring Ultra Search in a Hosted Environment	2-30
Preconfiguration Tasks for a Hosted Environment.....	2-30
Configuring Ultra Search in the Subscriber Context.....	2-30

3 Post-Installation Information

Changing Ultra Search Schema Passwords	3-2
Configuring the Oracle Server for Ultra Search	3-2
Step 1: Tune the Oracle Database.....	3-2
Step 2: Create and Assign the Temporary Tablespace to the CTXSYS User.....	3-4
Step 3: Create a Large Tablespace for Each Ultra Search Instance User.....	3-4
Step 4: Create and Configure New Database Users for Each Ultra Search Instance	3-5
Step 5: Alter the Index Preferences	3-6
Managing Stoplists.....	3-7
Default Ultra Search Stoplist.....	3-7
Modifying Instance Stoplists.....	3-7
Modifying Instance Stoplists Before Initial Crawling.....	3-8
Modifying Instance Stoplists After Initial Crawling	3-9
Upgrading Ultra Search	3-9
Pre-Upgrade Steps.....	3-10
Upgrading Ultra Search Shipped with Oracle Database	3-10
Upgrading Ultra Search Shipped with Oracle Application Server	3-10

Upgrading Ultra Search Shipped with Oracle Collaboration Suite	3-10
Upgrading Ultra Search to Oracle Collaboration Suite Release 1	3-11
Upgrade from Ultra Search 1.0.3 to 9.0.3	3-11
Upgrade from Ultra Search 9.0.2 to 9.0.3	3-14
Upgrade from Ultra Search 9.2 to 9.0.3	3-15
Configuring the Query Application	3-15
Step 1: Edit the data-sources.xml File	3-15
Step 2: Deploy Multiple Query Applications Against Multiple Instances.....	3-15

4 Tuning and Performance

Tuning the Web Crawling Process	4-2
Web Crawling Strategy	4-2
Monitoring the Crawling Process.....	4-2
URL Looping	4-2
Tuning Query Performance	4-3
Using the Remote Crawler	4-6
Scalability and Load Balancing.....	4-7
Installation and Configuration Sequence	4-7
Ultra Search on Real Application Clusters	4-10
Configuring Storage Access	4-10
Remote Crawler File Cache	4-11
Logging on to the Oracle Instance	4-12
Query Search Application for Read Application Clusters	4-12
Java Crawler	4-12
Choosing a JDBC Driver.....	4-12
Table Data Source Synchronization	4-13
Synchronizing Crawling of Oracle Databases.....	4-14
Create Log Table.....	4-14
Create Log Triggers.....	4-15
Synchronizing Crawling of Non-Oracle Databases.....	4-16

5 Security in Ultra Search

About Ultra Search Security	5-2
Ultra Search Security Model	5-2
Classes of Users and Their Privileges	5-3

Ultra Search Default Users	5-3
Ultra Search Admin Privilege Model in the Hosted Environment	5-4
Admin Privilege Model.....	5-5
Resources Protected by Ultra Search	5-7
Authorization and Access Enforcement.....	5-7
How Ultra Search Leverages Security Services	5-8
How Ultra Search Leverages the Identity Management Infrastructure.....	5-8
Ultra Search Extensibility and Security.....	5-9
Configuring a Security Framework for Ultra Search	5-9
Configuring Security Framework Options for Ultra Search.....	5-9
Configuring Oracle Identity Management Options for Ultra Search	5-9
Configuring Ultra Search Security	5-9

6 Understanding the Ultra Search Crawler and Data Sources

Overview of the Ultra Search Crawler	6-2
Crawler Settings	6-2
Crawler Data Sources	6-2
Using Crawler Agents.....	6-3
Synchronizing Data Sources	6-3
Display URL and Access URL	6-3
Document Attributes	6-3
Crawling Process for the Schedule	6-4
Queuing and Caching Documents.....	6-4
Indexing Documents.....	6-7
Data Synchronization	6-8
Ultra Search Remote Crawler	6-9

7 Understanding the Ultra Search Administration Tool

Ultra Search Administration Tool	7-1
Setting Crawler Parameters	7-2
Setting Query Options	7-3
Attributes.....	7-3
Data Groups.....	7-3
Online Help in Different Languages.....	7-3
Logging On to Ultra Search	7-3

Logging On and Managing Instances as SSO Users	7-5
Logging On to Ultra Search.....	7-5
Granting Privileges to SSO Users.....	7-5
Instances Page	7-6
Creating an Instance.....	7-7
Creating a Regular Instance.....	7-7
Creating a Snapshot Instance.....	7-8
Selecting an Instance.....	7-10
Deleting an Instance.....	7-11
Editing an Instance.....	7-11
Instance Mode.....	7-11
Schema Password.....	7-11
Crawler Page	7-12
Configure the Settings.....	7-12
Remote Crawler Profiles.....	7-16
Crawler Statistics.....	7-16
Summary of Crawler Activity.....	7-16
Detailed Crawler Statistics.....	7-16
Crawler Progress.....	7-16
Problematic URLs.....	7-16
Web Access Page	7-17
Proxies.....	7-17
Authentication.....	7-17
HTTP Authentication.....	7-17
HTML Forms.....	7-17
Attributes Page	7-18
Search Attributes.....	7-18
Mappings.....	7-19
Sources Page	7-20
Web Sources.....	7-20
Creating Web Sources.....	7-21
Table Sources.....	7-23
Creating Table Sources.....	7-24
Editing Table Sources.....	7-25
Table Sources Comprised of More Than One Table.....	7-25

Limitations With Database Links.....	7-25
Email Sources	7-26
Creating Email Sources.....	7-26
File Sources	7-27
Creating File Sources	7-27
Oracle Sources.....	7-28
Oracle Portal Sources.....	7-28
Federated Sources	7-29
User-Defined Sources.....	7-31
Creating User-Defined Data Source Types.....	7-31
Creating User-Defined Sources.....	7-32
Schedules Page	7-32
Data Synchronization.....	7-33
Creating Synchronization Schedules.....	7-33
Updating Schedules	7-33
Editing Synchronization Schedules.....	7-33
Launching Synchronization Schedules	7-35
Synchronization Status and Crawler Progress.....	7-36
Index Optimization	7-37
Queries Page	7-38
Data Groups	7-38
URL Submission	7-38
Relevancy Boosting	7-39
Query Statistics	7-40
Configuration	7-41
Users Page	7-42
Preferences.....	7-42
Super-Users	7-42
Privileges.....	7-43
Globalization Page	7-43
Search Attribute Name	7-43
LOV Display Name	7-45
Data Group Name	7-45

8 Ultra Search Developer's Guide and API Reference

Overview of Ultra Search APIs	8-2
Ultra Search Query API	8-2
Customizing the Query Syntax Expansion	8-3
Default Query Syntax Expansion Implementation.....	8-4
End User Query Syntax	8-4
Scoring Classes.....	8-6
Expansion Rules.....	8-7
Examples of Applying the Rules.....	8-7
Customizing the Rules	8-8
Ultra Search Query Tag Library	8-9
Query Tag Descriptions	8-11
<instance> Tag: Connecting to the Ultra Search Instance	8-11
<iterAttributes> Tag: Show All Search Attributes	8-13
<iterGroups> Tag: Show All Search Groups	8-13
<iterLanguages> Tag: Show All Search Languages	8-14
<iterLOV> Tag: Show All Values Defined for a Search Attribute	8-15
Formulating the Query	8-15
<getResult> Tag: Perform Search.....	8-15
<fetchAttribute> Tag: Metadata Selection	8-16
<showHitCount> Tag: Show Estimated Hit Count.....	8-17
<iterResult> Tag: Render the Results	8-18
<showAttributeValue> Tag: Render a Document Attribute.....	8-18
Ultra Search Crawler Agent API	8-19
Crawler Agent Overview	8-19
Standard Agent	8-20
Smart Agent.....	8-21
Document Attributes and Properties.....	8-21
Crawler Agent Functionality	8-21
Data Source Type Registration	8-21
Data Source Registration	8-22
Data Source Attribute Registration.....	8-22
User-Implemented Crawler Agent	8-23
Interaction Between the Crawler and the Crawler Agent	8-23
Crawler Agent APIs and Classes	8-23

Sample Agent Files.....	8-24
Setting up the Sample Crawler Agent.....	8-24
Compiling and Building the Agent Jar File.....	8-24
Creating a Data Source Type.....	8-25
Defining Data Source Parameters.....	8-25
Defining a Data Source of this Type.....	8-25
Ultra Search Java Email API.....	8-26
JavaMail Implementation.....	8-27
Java Email API.....	8-27
Sample Mailing List Browser Application Files.....	8-28
Setting up the Sample Mailing List Browser Application.....	8-29
Ultra Search URL Rewriter API.....	8-29
URL Link Filtering.....	8-29
URL Link Rewriting.....	8-30
Creating and Using a URL Rewriter.....	8-32
Ultra Search Sample Query Applications.....	8-33
Sample Query Applications.....	8-34
JavaServer Page Concepts.....	8-34

A Loading Metadata into Ultra Search

Launching the Loading Tool.....	A-1
Loading Documents and Relevance Scores.....	A-2
The Input XML File.....	A-2
Example of the Document Relevance Boosting XML File.....	A-3
Loading Search Attribute LOVs and LOV Display Names.....	A-3
The LOV XML File.....	A-3
Example of the LOV XML File.....	A-4
XML Schema for Document Relevance Boosting.....	A-5
XML Schema for LOVs and LOV Display Names.....	A-5

B Altering the Crawler Java Classpath

Reasons for Altering the Crawler Java Classpath.....	B-1
Difference Between the Crawler Classpath and the Remote Crawler Classpath.....	B-1
Altering the Crawler Java Classpath on the Ultra Search Server Host.....	B-2
Altering the Crawler Java Classpath on a Remote Crawler Host.....	B-2

C Customizing the Query Syntax Expansion 9.0.1

Default Query Syntax Expansion Implementation	C-1
End User Query Syntax	C-1
Summary of Rules	C-2
Scoring.....	C-3
Expansion Rules.....	C-4
Customizing the Rules	C-5
The expand_main Function	C-6
The expand_attr Function	C-6
Example of Combining Values.....	C-7

Index

Send Us Your Comments

Oracle Ultra Search User's Guide 10g (9.0.4)

Part No. B10896-01

Oracle Corporation welcomes your comments and suggestions on the quality and usefulness of this publication. 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 about this manual?

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

- Electronic mail: infodev_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 electronic mail address (optional).

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

Preface

Oracle Ultra Search User's Guide describes how to configure and use Ultra Search and Ultra Search APIs.

This preface contains these topics:

- Audience
- Organization
- Related Documentation
- Conventions
- Documentation Accessibility

Audience

Oracle Ultra Search User's Guide is intended for database administrators and application developers who perform the following tasks:

- Install and configure Ultra Search
- Administer Oracle Ultra Search instances
- Develop Oracle Ultra Search applications

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

Organization

This document contains:

"What's New in Ultra Search?"

This section describes new features and provides pointers to additional information.

Chapter 1, "Introduction to Ultra Search"

This chapter provides an overview of Ultra Search and describes the system configuration.

Chapter 2, "Installing and Configuring Ultra Search"

This chapter describes how to install and configure Ultra Search.

Chapter 5, "Security in Ultra Search"

This chapter describes the architecture and configuration of security for Ultra Search.

Chapter 3, "Post-Installation Information"

This chapter provides post-installation information, such as how to configure the Oracle server for Ultra Search and how to manage stoplists. It also describes how to upgrade to the most recent Ultra Search release.

Chapter 4, "Tuning and Performance"

This chapter describes various ways to tune Ultra Search and improve performance. These include tuning the Web crawling process, tuning query performance, using the remote crawler, using Ultra Search on Real Application Clusters, and table data source synchronization.

Chapter 6, "Understanding the Ultra Search Crawler and Data Sources"

This chapter explains how the crawler works. It also describes crawler settings, data sources, document attributes, data synchronization, and the remote crawler.

Chapter 7, "Understanding the Ultra Search Administration Tool"

This chapter describes how to use the Ultra Search administration tool to configure and schedule the Ultra Search crawler.

Chapter 8, "Ultra Search Developer's Guide and API Reference"

This chapter explains the following Ultra Search APIs: query API, crawler agent API, email API, and URL rewriter API. It also provides related API information, such as details about the sample query applications, the query tag library, and query syntax expansion customization.

Appendix A, "Loading Metadata into Ultra Search"

This appendix describes the command-line tool for loading metadata into an Ultra Search database.

Appendix B, "Altering the Crawler Java Classpath"

This appendix explains why and how to alter the crawler Java classpath.

Appendix C, "Customizing the Query Syntax Expansion 9.0.1"

This appendix describes how to customize the query syntax expansion for Ultra Search release 1 (9.0.1).

Related Documentation

For more information, see these Oracle resources:

- *Oracle9i Database Concepts*
- *Oracle9i Database Administrator's Guide*
- *Oracle9i Database Performance Tuning Guide and Reference*
- *Oracle Enterprise Manager Concepts*

Many books in the documentation set use the sample schemas of the seed database, which is installed by default when you install Oracle. Refer to *Oracle9i Sample Schemas* for information on how these schemas were created and how you can use them yourself.

Printed documentation is available for sale in the Oracle Store at

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

To download free release notes, installation documentation, 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://otn.oracle.com/membership/>

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

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

To access the database documentation search engine directly, please visit

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

Conventions

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

- Conventions in Text
- Conventions in Code Examples
- Conventions for Windows Operating Systems

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.	When you specify this clause, you create an index-organized table .
<i>Italics</i>	Italic typeface indicates book titles or emphasis.	<i>Oracle9i Database Concepts</i> Ensure that the recovery catalog and target database do <i>not</i> reside on the same disk.
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.	You can specify this clause only for a NUMBER column. You can back up the database by using the BACKUP command. Query the TABLE_NAME column in the USER_TABLES data dictionary view. Use the DBMS_STATS.GENERATE_STATS procedure.

Convention	Meaning	Example
lowercase monospace (fixed-width) font	Lowercase monospace typeface indicates executables, filenames, directory names, 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. Note: Some programmatic elements use a mixture of UPPERCASE and lowercase. Enter these elements as shown.	Enter <code>sqlplus</code> to open SQL*Plus. The password is specified in the <code>orapwd</code> file. Back up the datafiles and control files in the <code>/disk1/oracle/dbs</code> directory. The <code>department_id</code> , <code>department_name</code> , and <code>location_id</code> columns are in the <code>hr.departments</code> table. Set the <code>QUERY_REWRITE_ENABLED</code> initialization parameter to <code>true</code> . Connect as <code>oe</code> user. The <code>JRepUtil</code> class implements these methods.
lowercase italic monospace (fixed-width) font	Lowercase italic monospace font represents placeholders or variables.	You can specify the <code>parallel_clause</code> . Run <code>old_release.SQL</code> where <code>old_release</code> refers to the release you installed prior to upgrading.

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 (<i>digits</i> [, <i>precision</i>])
{ }	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.	<pre>SQL> SELECT NAME FROM V\$DATAFILE; NAME ----- /fs1/dbs/tbs_01.dbf /fs1/dbs/tbs_02.dbf . . . /fs1/dbs/tbs_09.dbf 9 rows selected.</pre>
Other notation	You must enter symbols other than brackets, braces, vertical bars, and ellipsis points as shown.	<pre>acctbal NUMBER(11,2); acct CONSTANT NUMBER(4) := 3;</pre>
<i>Italics</i>	Italicized text indicates placeholders or variables for which you must supply particular values.	<pre>CONNECT SYSTEM/system_password DB_NAME = database_name</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. Note: Some programmatic elements use a mixture of UPPERCASE and lowercase. Enter these elements as shown.	<pre>SELECT last_name, employee_id FROM employees; sqlplus hr/hr CREATE USER mjones IDENTIFIED BY ty3MU9;</pre>

Conventions for Windows Operating Systems

The following table describes conventions for Windows operating systems and provides examples of their use.

Convention	Meaning	Example
Choose Start >	How to start a program.	To start the Database Configuration Assistant, choose Start > Programs > Oracle - <i>HOME_NAME</i> > Configuration and Migration Tools > Database Configuration Assistant.
File and directory names	File and directory names are not case sensitive. The following special characters are not allowed: left angle bracket (<), right angle bracket (>), colon (:), double quotation marks ("), slash (/), pipe (), and dash (-). The special character backslash (\) is treated as an element separator, even when it appears in quotes. If the file name begins with \\, then Windows assumes it uses the Universal Naming Convention.	c:\winnt"\system32 is the same as C:\WINNT\SYSTEM32
C:\>	Represents the Windows command prompt of the current hard disk drive. The escape character in a command prompt is the caret (^). Your prompt reflects the subdirectory in which you are working. Referred to as the <i>command prompt</i> in this manual.	C:\oracle\oradata>
Special characters	The backslash (\) special character is sometimes required as an escape character for the double quotation mark (") special character at the Windows command prompt. Parentheses and the single quotation mark (') do not require an escape character. Refer to your Windows operating system documentation for more information on escape and special characters.	C:\>exp scott/tiger TABLES=emp QUERY=\"WHERE job='SALESMAN' and sal<1600\" C:\>imp SYSTEM/password FROMUSER=scott TABLES=(emp, dept)
<i>HOME_NAME</i>	Represents the Oracle home name. The home name can be up to 16 alphanumeric characters. The only special character allowed in the home name is the underscore.	C:\> net start Oracle <i>HOME_NAME</i> TNSListener

Convention	Meaning	Example
<i>ORACLE_HOME</i> and <i>ORACLE_BASE</i>	<p>In releases prior to Oracle8i release 8.1.3, when you installed Oracle components, all subdirectories were located under a top level <i>ORACLE_HOME</i> directory that by default used one of the following names:</p> <ul style="list-style-type: none"> ■ C:\orant for Windows NT ■ C:\orawin98 for Windows 98 <p>This release complies with Optimal Flexible Architecture (OFA) guidelines. All subdirectories are not under a top level <i>ORACLE_HOME</i> directory. There is a top level directory called <i>ORACLE_BASE</i> that by default is C:\oracle. If you install the latest Oracle release on a computer with no other Oracle software installed, then the default setting for the first Oracle home directory is C:\oracle\orann, where <i>nn</i> is the latest release number. The Oracle home directory is located directly under <i>ORACLE_BASE</i>.</p> <p>All directory path examples in this guide follow OFA conventions.</p> <p>Refer to <i>Oracle9i Database Getting Started for Windows</i> for additional information about OFA compliances and for information about installing Oracle products in non-OFA compliant directories.</p>	Go to the <i>ORACLE_BASE\ORACLE_HOME\rdbms\admin</i> directory.

Documentation Accessibility

Our goal is to make Oracle products, services, and supporting documentation accessible, with good usability, to the disabled community. 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 Corporation 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/>

Accessibility of Code Examples in Documentation

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.

Accessibility of Links to External Web Sites in Documentation

This documentation may contain links to Web sites of other companies or organizations that Oracle Corporation does not own or control. Oracle Corporation neither evaluates nor makes any representations regarding the accessibility of these Web sites.

What's New in Ultra Search?

This section describes Ultra Search new features, with pointers to additional information. It also explains the Ultra Search release history.

Secure Crawling

Ultra Search provides secure crawling with the following types of authentication:

Digest Authentication Ultra Search supports HTTP digest authentication, and the Ultra Search crawler can authenticate itself to Web servers employing HTTP digest authentication scheme. This is based on a simple challenge-response paradigm; however, the password is encrypted.

HTML Form Authentication HTML form-based authentication is the most commonly used authentication scheme on the Web. Ultra Search lets you register HTML forms that you want the Ultra Search crawler to automatically fill out during Web crawling. HTML form authentication requires that HTTP cookie functionality is enabled.

See Also: "Creating Web Sources" on page 7-21

Indexing Dynamic Pages

Dynamic URLs can be crawled and indexed. Some dynamic pages appear as multiple search hits for the same page, and you may not want them all indexed. Other dynamic pages are each different and need to be indexed.

See Also: "Creating Web Sources" on page 7-21

HTTPS

Ultra Search now supports HTTPS (HTTP over SSL). The Ultra Search crawler can now crawl HTTPS URLs (for example, <https://www.foo.com>).

See Also: "Creating Web Sources" on page 7-21

Secure Searching

Ultra Search now supports secure searches. Secure searches return only documents that the search user is allowed to view.

Each indexed document can be protected by an access control list (ACL). During searches, the ACL is evaluated. If the user performing the search has permission to read the protected document, then the document is returned by the query API. Otherwise, it is not returned.

Ultra Search stores ACLs in the Oracle XML DB repository. Ultra Search also uses Oracle XML DB functionality to evaluate ACLs.

See Also: "Secure Search" on page 1-6

Integration with Oracle Internet Directory

Oracle Internet Directory (OID) is Oracle's native LDAP v3-compliant directory service, built as an application on top of the Oracle database. Ultra Search integrates with OID in the following areas:

- Ultra Search administration groups and group membership are stored in OID.
- Users are authenticated through the single sign-on (SSO) server and OID.
- OID performs authorization on Ultra Search users' administration privileges.

See Also: "Integration with Oracle Internet Directory" on page 1-11

Cookie Support

Cookie support is enabled by default.

Crawler Cache Deletion Control

During crawling, documents are stored in the cache directory. Every time the preset size is reached, crawling stops and indexing starts. In previous releases, the cache file was always deleted when indexing was done. You can now specify *not* to delete

the cache file when indexing is done. This option applies to all data sources. The default is to delete the cache file after indexing.

See Also: "Crawler Page" on page 7-12

URL Boundary Rules Include Port Number Inclusion or Exclusion

You can set URL boundary rules to refine the crawling space. You can now include or exclude Web sites with a specific port. For example, you can include `www.oracle.com` but not `www.oracle.com:8080`. By default, all ports are crawled.

See Also: "Creating Web Sources" on page 7-21

Hostname Prefix Allowed in Web Data Source URL Boundary Specification

In previous releases, you could only specify suffix inclusion rules. For example, crawl only URLs ending with `"oracle.com"`. You can now also specify prefix rules. For example, crawl `"oracle.com"` but not `"stores.oracle.com"`.

See Also: "Creating Web Sources" on page 7-21

Default Ultra Search Instance and Schema

Ultra Search automatically creates a default Ultra Search instance based on the default Ultra Search test user. So, you can test Ultra Search functionality based on the default instance after installation.

See Also: "Configuring the Default Ultra Search Instance" on page 2-11

Crawler Recrawl Policy

You can update the recrawl policy to process documents that have changed or to process all documents.

In previous releases, "process all documents" did not help when the crawling scope had been narrowed. For example, if crawling depth was reduced from seven to five, the PDF mimetype was deleted, or a host inclusion rule was removed, then you had to remove the affected documents manually in a SQL*Plus session.

With this release, all crawled URLs are subject to crawler setting enforcement, not just newly crawled URLs.

See Also: "Editing Synchronization Schedules" on page 7-33

Federated Search

Traditionally, Oracle Ultra Search used *centralized* search to gather data on a regular basis and update one index that cataloged all searchable data. This provided fast searching, but it required that the data source to be crawlable before it could be searched. Ultra Search now also provides *federated* search, which allows multiple indexes to perform a single search. Each index can be maintained separately. By querying the data source at search-time, search results are always the latest results. User credentials can be passed to the data source and authenticated by the data source itself. Queries can be processed efficiently using the data's native format.

To use federated search, you must deploy an Ultra Search search adapter, or *searchlet*, and create an Oracle source. A searchlet is a Java module deployed in the middle tier (inside OC4J) that searches the data in an enterprise information system on behalf of a user. When a user's query is delegated to the searchlet, the searchlet runs the query on behalf of the user. Every searchlet is a JCA 1.0 compliant resource adapter.

See Also: "Federated Sources" on page 7-29

Ultra Search Release Information

Ultra Search is released with the Oracle Database, Oracle Application Server, and Oracle Collaboration Suite. Because of different release numbers in the past, the Ultra Search release numbers are somewhat confusing.

- Ultra Search 9.0.4 is part of Oracle Application Server (OracleAS) release 10g (9.0.4).
- Ultra Search release 9.0.3 is part of the Oracle Collaboration Suite release 9.0.3.
- Ultra Search release 9.2 is part of Oracle9i release 9.2. Ultra Search release 1.0.3 was part of Oracle9i release 1 (9.0.1).
- Ultra Search release 9.0.2 is part of Oracle9iAS release 2 (9.0.2).

Introduction to Ultra Search

This chapter contains the following topics:

- Overview of Ultra Search
- Ultra Search Components
- Ultra Search Features
- Ultra Search System Configuration

Overview of Ultra Search

Oracle Ultra Search is built on the Oracle Database and Oracle Text technology that provides uniform search-and-locate capabilities over multiple repositories: Oracle databases, other ODBC compliant databases, IMAP mail servers, HTML documents served up by a Web server, files on disk, and more.

Ultra Search uses a 'crawler' to collect documents. You can schedule the crawler to suit the Web sites that you want to search. The documents stay in their own repositories, and the crawled information is used to build an index that stays within your firewall in a designated Oracle database. Ultra Search also provides APIs for building content management solutions.

In addition, Ultra Search offers the following:

- A complete text query language for text search inside the database
- Full integration with the Oracle Database and the SQL query language
- Advanced features like concept searching and theme analysis
- Attribute mapping to facilitate attribute search across disparate repositories
- Indexing of all popular file formats (150+)
- Full globalization, including support for Chinese, Japanese and Korean (CJK), and Unicode

Ultra Search Components

Ultra Search is made up of the following components:

- Ultra Search Crawler
- Ultra Search Backend
- Ultra Search Administration Tool
- Ultra Search APIs and Sample Applications

Ultra Search Crawler

The Ultra Search crawler is a Java process activated by your Oracle server according to a set schedule. When activated, the crawler spawns a configurable number of processor threads that fetch documents from various data sources and index them using Oracle Text. This index is used for querying. Data sources can be Web sites,

database tables, files, mailing lists, OracleAS Portal page groups, or user-defined data sources.

The crawler maps links and analyzes relationships. The crawler schedule is integrated with and driven from the `DBMS_JOB` queue mechanism. Whenever the crawler encounters embedded, non-HTML documents during the crawling, it uses Oracle Text filters to automatically detect the document type and filter and index the document.

See Also: Chapter 6, "Understanding the Ultra Search Crawler and Data Sources"

Ultra Search Backend

The Ultra Search backend (server component) consists of an Ultra Search repository and Oracle Text. Oracle Text provides the text indexing and search capabilities required to index and query data retrieved from your data sources. The backend is not visible to users; it indexes information from the crawler and serves up the query results.

See Also: "Installing the Ultra Search Backend" on page 2-4

Ultra Search Administration Tool

The administration tool is a J2EE-compliant Web application. You can use it to manage Ultra Search instances, and you can access it from any browser in your intranet. The administration tool is independent from the Ultra Search query application. Therefore, the administration tool and query application can be hosted on different computers to enhance security and scalability.

See Also: Chapter 7, "Understanding the Ultra Search Administration Tool"

Ultra Search APIs and Sample Applications

Ultra Search provides the following APIs:

- The query API works with indexed data. The Java API does not impose any HTML rendering elements. The application can completely customize the HTML interface.
- The crawler agent API crawls and indexes proprietary document repositories.

- The email Java API accesses archived emails and is used by the query application to display emails. It can also be used when building your own custom query application.
- The URL rewriter API is used by the crawler to filter and rewrite extracted URL links before they are inserted into the URL queue.

Ultra Search includes highly functional query applications to query and display search results. The query applications are based on JSP and work with any JSP1.1 compliant engine.

See Also:

- Chapter 8, "Ultra Search Developer's Guide and API Reference"
- *Oracle Ultra Search API Reference*

Ultra Search Features

This section explains some features in Ultra Search. It includes the following topics:

- Integration with Oracle Application Server
- Extensible Crawler and Crawler Agents
- Federated Search
- Secure Search
- Sample Query Applications
- Sample Search Portlet
- Query API
- URL Rewrite
- Robots Exclusions
- Display URL Support
- Document and Search Attributes
- Metadata Loader
- Document Relevancy Boosting
- Data Harvesting Mode
- Instance Snapshot Support

- Integration with Oracle Internet Directory
- Single Sign-On Authentication
- Query Syntax Expansion

Integration with Oracle Application Server

Although Ultra Search in the Oracle Application Server (OracleAS) is the same product as Ultra Search in Oracle Collaboration Suite and Ultra Search in the Oracle Database, there are a couple differences:

- The Oracle Database is not integrated with OracleAS Portal. With OracleAS and Oracle Collaboration Suite, Portal users add powerful multi-repository search to their Portal pages. OracleAS and Oracle Collaboration Suite also have the capability to crawl and make searchable Portal's own repository.
- OracleAS includes a Single Sign-On (SSO) server. SSO users can log on once for all components of the OracleAS product, and the Ultra Search administrative interface allows user management operations on either database users or SSO users. Authenticated SSO users never see the Ultra Search logon screen. Instead, they can immediately choose an instance. If the SSO user does not have permissions to manage Ultra Search (set in the **Users Page**), then the SSO user receives an error. SSO is available only with the Oracle Identity Management infrastructure.

See Also: <http://portalstudio.oracle.com/>

Extensible Crawler and Crawler Agents

You can define, edit, or delete your own data sources and types in addition to the ones provided. You might implement your own crawler agent to crawl and index a proprietary document repository, such as Lotus Notes or Documentum, which contain their own databases and interfaces. The proprietary repository is called a user-defined data source. The module that enables the crawler to access the data source is called a crawler agent.

See Also:

- "Ultra Search Crawler Agent API" on page 8-19
- *Oracle Ultra Search API Reference*

Federated Search

Traditionally, Oracle Ultra Search used *centralized* search to gather data on a regular basis and update one index that cataloged all searchable data. This provided fast searching, but it required that the data source to be crawlable before it could be searched. Ultra Search now also provides *federated* search, which allows multiple indexes to perform a single search. Each index can be maintained separately. By querying the data source at search-time, search results are always the latest results. User credentials can be passed to the data source and authenticated by the data source itself. Queries can be processed efficiently using the data's native format.

To use federated search, you must deploy an Ultra Search search adapter, or *searchlet*, and create an Oracle source. A searchlet is a Java module deployed in the middle tier (inside OC4J) that searches the data in an enterprise information system on behalf of a user. When a user's query is delegated to the searchlet, the searchlet runs the query on behalf of the user. Every searchlet is a JCA 1.0 compliant resource adapter.

See Also: "Federated Sources" on page 7-29

Secure Search

Ultra Search supports secure searches, which return only documents satisfying the search criteria that the search user is allowed to view. For secure searches, each indexed document should be protected by an access control list (ACL). During searches, the ACL is evaluated. If the user performing the search has permission to read the protected document, then the document is returned by the query API. Otherwise, it is not returned.

There are two ways to secure a data source:

- Specify a single ACL for protecting all documents of a data source.
The administrator specifies the permissions of the single ACL in the Ultra Search administration tool. The resulting ACL is used to protect all documents belonging to that data source.
- Crawl ACLs from the data source.
The data source is expected to provide the ACL together with the document. This lets each document be protected by its own unique ACL.

Ultra Search performs ACL duplicate detection. This means that if a crawled document's ACL already exists in the Ultra Search system, then the existing ACL is used to protect the document, instead of creating a new ACL within Ultra Search. This policy reduces storage space and increases performance.

Ultra Search supports only a single LDAP domain. The LDAP users and groups specified in the ACL must belong to the same LDAP domain.

Caution: If ACLs are crawled from data sources, then it is the responsibility of the administrator to ensure that the data sources being crawled belong to the same LDAP domain. Otherwise, it is possible that search users can inadvertently be granted permissions to access documents that they should not be able to access.

Searches executed against a secure-search enabled Ultra Search instance are slower than those executed against a non secure-search enabled instance. This is because each candidate hit could require an ACL evaluation. ACLs are evaluated natively by the Oracle server for optimum performance. Nevertheless, this is a finite time. Therefore, the time taken to return hits in a secure search varies depending on the number ACL evaluations that must be made.

Dependency on Oracle XML DB

Ultra Search stores ACLs in the Oracle XML DB repository. Ultra Search also uses Oracle XML DB functionality to evaluate ACLs.

The ACLs are managed by Ultra Search. ACLs are uniquely referenced by documents from a single Ultra Search instance. ACLs are not shared by multiple Ultra Search instances. For acceptable performance, the ACL cache size must be large enough to contain all ACLs evaluated at run time.

ACLs in the XML DB repository are protected by other ACLs (known as "protector ACLs"). Ultra Search ensures that the protector ACLs grant appropriate privileges in order for Ultra Search to invoke the XML DB ACL evaluation mechanism. The evaluation performance is primarily affected by the total number of ACLs used by all XML DB client applications that also utilize its ACL evaluation mechanism. This set of applications includes Ultra Search.

See Also: *Oracle9i XML Database Developer's Guide - Oracle XML DB* and "Configure a Secure Ultra Search Installation" on page 2-6

Sample Query Applications

Ultra Search includes fully functional sample query applications to query and display search results. The sample query applications include a sample search portlet. The sample Ultra Search portlet demonstrates how to write a search portlet

for use in OracleAS Portal. This same portlet is installed as a feature of the OracleAS Portal product.

See Also: "Ultra Search Query API" on page 8-2

Sample Search Portlet

Ultra Search provides a search portlet that can be embedded in OracleAS Portal pages. It is implemented as a JavaServer Page application.

The Ultra Search search portlet supports most of the functionality provided by the Query API Complete Sample application.

See Also:

- The OracleAS Portal documentation for more information about portlets
- Oracle Ultra Search Sample Query Applications Readme for more information about the Query API Complete Sample application

Query API

Oracle Ultra Search offers a flexible API to incorporate search functionality to your portal site. The query API includes the following functionality:

- Three attribute types: string, number, and date
- Multivalued attributes
- Display name support for attributes, attribute list of values (LOV), and data groups
- Document relevancy boosting
- Arbitrary grouping of attribute query operator using operators (AND, OR), with control over attribute operator evaluation order
- Selection of metadata returned in query result

See Also:

- "Ultra Search Query API" on page 8-2
- *Oracle Ultra Search API Reference*

URL Rewrite

The URL rewriter is a user-supplied Java module for implementing the Ultra Search `UrlRewriter` interface. It is used by the crawler to filter or rewrite extracted URL links before they are put into the URL queue. URL filtering removes unwanted links, and URL rewriting transforms the URL link. This transformation is necessary when access URLs are used.

See Also:

- "Web Sources" on page 7-20
- "Ultra Search URL Rewriter API" on page 8-29
- *Oracle Ultra Search API Reference*

Robots Exclusions

Robots exclusion lets you control which parts of your sites can be visited by robots. If robots exclusion is enabled (default), then the Web crawler traverses the pages based on the access policy specified in the Web server `robots.txt` file. For example, when a robot visits `http://www.foobar.com/`, it checks for `http://www.foobar.com/robots.txt`. If it finds it, the crawler analyzes its contents to see if it is allowed to retrieve the document. If you own the Web sites, then you can disable robots exclusions. However, when crawling other Web sites, you should always comply with `robots.txt` by enabling robots exclusion.

See Also: "Web Sources" on page 7-20

Display URL Support

When gathering information from a database-based Web application, Ultra Search lets you specify a URL to display the data retrieved on a browser, rendered by a screen of a Web application corresponding to the data in the database tables. The URL points to a screen in the Web application corresponding to the data in the database. This is available for table data sources, file data sources, and user-defined data sources.

See Also: "Using Crawler Agents" on page 6-3

Document and Search Attributes

Document attributes, or metadata, describe the properties of a document. Each data source has its own set of document attributes. The value is retrieved during the

crawling process and then mapped to one of the search attributes and stored and indexed in the database. This lets you query documents based on their attributes. Document attributes in different data sources can be mapped to the same search attribute. Therefore, you can query documents from multiple data sources based on the same search attribute.

The list of values (LOV) for a search attribute can help you specify a search query. If attribute LOV is available, then the crawler registers the LOV definition, which includes attribute value, attribute value display name, and its translation.

See Also: "Synchronizing Data Sources" on page 6-3

Metadata Loader

Ultra Search provides a command-line tool to load metadata into an Ultra Search database. If you have a large amount of data, this is probably faster than using the HTML-based administration tool. The loader tool supports the following types of metadata:

- Search attribute list of values (LOVs) and display names
- Document relevance boosting and document loading

See Also: Appendix A, "Loading Metadata into Ultra Search"

Document Relevancy Boosting

You can override the search results and influence the order that documents are ranked in the query result list with document relevancy boosting. This can promote important documents to higher scores and make them easier to find.

Relevancy boosting assigns a score to a document for specific queries entered by the search user.

Note: The document still has a score computed by Oracle Text if you enter a query that is not one of the boosted queries.

Relevancy boosting has the following limitations:

- Comparison of the user's query against the boosted queries uses exact string match. This means that the comparison is case-sensitive and space-aware. Therefore, a document with a boosted score for "Ultra Search" is not boosted when you enter "ultrashow".

- Relevancy boosting requires that the query application pass in the search term in the Query API `getResult()` method call. The sample applications are designed to pass the basic search terms as the boost term. Advanced search criteria based on search attributes are ignored.

See Also: "Queries Page" on page 7-38

Data Harvesting Mode

For initial planning purposes, you might want the crawler to collect URLs without indexing. After crawling is done, you can examine document URLs and status, remove unwanted documents, and start indexing. You can update the crawling mode to the following:

- Automatically accept all URLs for indexing
- Examine URLs before indexing
- Index only

See Also: "Schedules Page" on page 7-32

Instance Snapshot Support

You can create a read-only snapshot of a master Ultra Search instance. This is useful for query processing or for a backup. You can also make a snapshot instance updatable. This is useful when the master instance is corrupted and you want to use a snapshot as a new master instance.

See Also: "Instances Page" on page 7-6

Integration with Oracle Internet Directory

Oracle Internet Directory (OID) is Oracle's native LDAP v3-compliant directory service, built as an application on top of the Oracle Database. OID hosts the Oracle common identity. All Oracle Web-based products integrate with the SSO server for single sign-on support.

Ultra Search Administration Groups in OID

An Ultra Search administration group contains a set of users. Each user can belong to one or multiple groups. All groups are created using `groupOfUniqueNames` and `orclGroup` object classes.

The only way to grant a user administration privileges is to assign them to an administration group. Ultra Search authorizes the user administration privileges based on the administration groups to which the user belongs. The following groups are created for each Ultra Search instance:

- **Super-users:** Users in this group can create or drop Ultra Search instances and can administer Ultra Search instances within the installation. Super-users must obey the rules for document relevancy boosting and ACL defined for each of the documents associated with the Ultra Search instance. For example, if a document ACL does not grant access to the super-user or group, then the super-user cannot search and browse the document.
- **Instance administrators:** Users in this group can administer the Ultra Search instance. Only the instance database schema user and members in the super-users group can drop the instance.

Authorization of the Administration Privileges

The authorization of the administration user is performed in the following steps:

1. After the administration user is successfully authenticated by the SSO server or the Ultra Search database, the Ultra Search GUI brings up the first screen for the user to choose an Ultra Search instance.
2. The Ultra Search GUI looks up the OID server or Ultra Search repository to find all Ultra Search instances with the installation that the administration user has privileges to administer.
3. The administration user chooses the Ultra Search instance from the list.

See Also: *Oracle Internet Directory Administrator's Guide*

Single Sign-On Authentication

The Ultra Search administration tool supports three modes of logging on, depending on the type of user. You can log on as:

- A single sign-on (SSO) user managed in the Oracle Internet Directory (OID) and authenticated with the SSO server
- A local database schema user in the Ultra Search database (non-SSO mode)
- A Portal user
- An Enterprise Manager user

Note: Single Sign-On (SSO) is available only with the Oracle Identity Management infrastructure.

See Also: "Logging On to Ultra Search" on page 7-3

Query Syntax Expansion

Ultra Search translates each user query into a database query. This process is called query syntax expansion. The expansion logic determines relevancy, recall of the search results. The Ultra Search default expansion boosts the relevancy of those documents that matches the user's query as a part of their title.

The query syntax expansion can be customized with the query API.

See Also: "Customizing the Query Syntax Expansion" on page 8-3

Ultra Search System Configuration

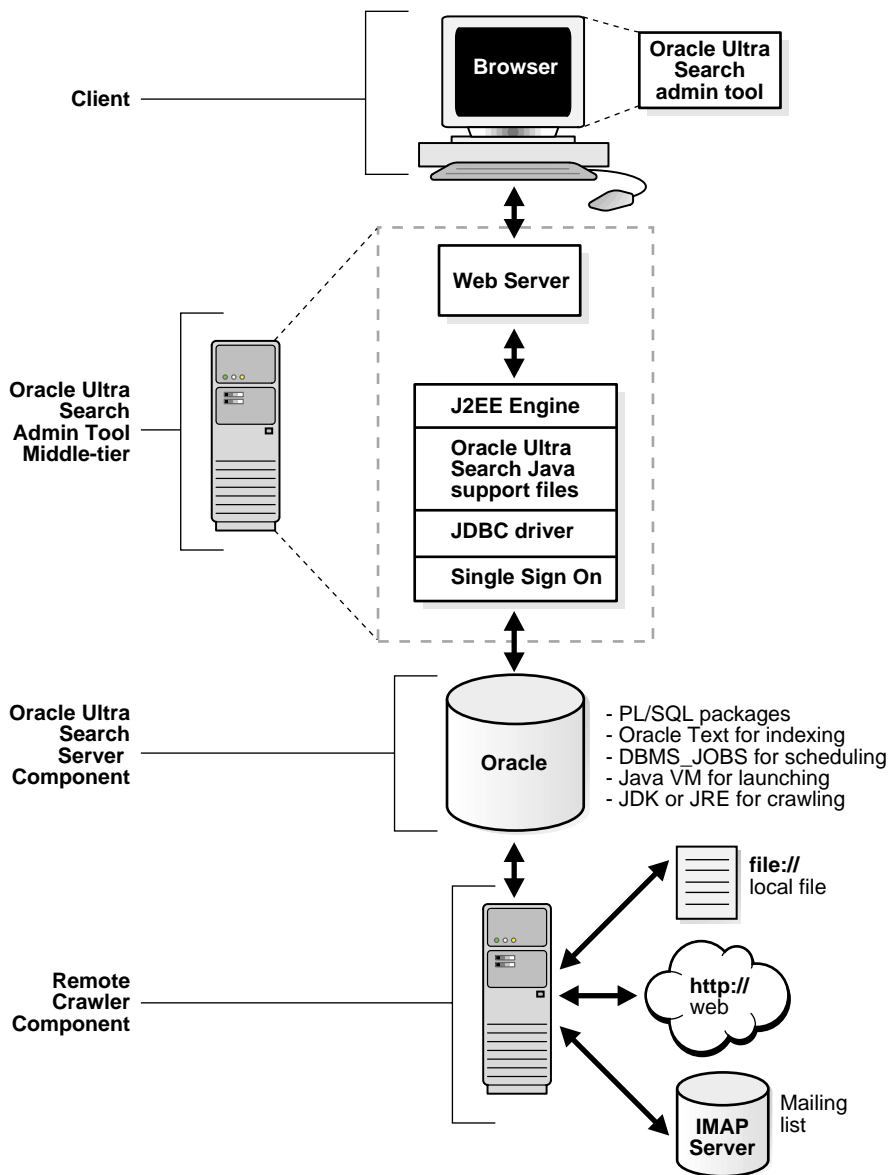
Ultra Search is a client program to the Oracle server at run time. It can be deployed in two configurations: in the backend or in the middle tier.

The Ultra Search query interface and the administration tool can be accessed from any HTML browser client. The administration tool relies on certain Java classes in the middle tier. This logical middle tier can be the same physical computer as the one that runs the database server, or a different one, running Oracle Application Server. The Ultra Search database backend (server component) consists of the Ultra Search data dictionary that stores metadata on all the different repositories, as well as the schedules and Java classes needed to drive the crawler. The crawler itself can run either on the database server computer or remotely on another computer.

See Also: Chapter 2, "Installing and Configuring Ultra Search" for more information about the components

Figure 1-1 illustrates the Ultra Search system configuration.

Figure 1-1 Oracle Ultra Search System Configuration



Installing and Configuring Ultra Search

This chapter contains the following topics:

- Ultra Search Requirements
- Installing the Ultra Search Backend
- Installing the Backend on an Existing Database or Metadata Repository
- Configuring the Default Ultra Search Instance
- Installing the Ultra Search Middle Tier on Web Server Hosts
- Installing the Backend on Remote Crawler Hosts
- Configuring Ultra Search in a Hosted Environment

Ultra Search Requirements

To use Oracle Ultra Search, you must install the following components with the Oracle Universal Installer:

- Ultra Search backend (server component)
 - Ultra Search repository
 - Ultra Search crawler
- Ultra Search middle tier
 - Ultra Search administration tool
 - Ultra Search query API

Note: By default, the Ultra Search crawler resides in the ORACLE_HOME of the database that contains the Ultra Search repository. You can install additional Ultra Search crawler components on other Oracle homes, usually on remote hosts. These are called remote crawlers, and they provide an effective way to deal with the scalability of the system.

The Ultra Search backend (server component) can be installed on any already existing database that is Oracle Database release 9.0.1.4 or higher with Oracle Text installed.

See Also: "Installing the Backend on an Existing Database or Metadata Repository" on page 2-9

Ultra Search Conventions

ORACLE_HOME refers to the Oracle home directory containing the Ultra Search backend (server component) or middle tier bits.

REMOTE_ORACLE_HOME refers to the Oracle home directory containing the Ultra Search remote crawler bits.

Ultra Search System Requirements

The following section describes the Ultra Search system requirements.

Hardware Requirements

Ultra Search hardware requirements vary based on the quantity of data that you plan to process using Ultra Search. Ultra Search uses Oracle Text as its indexing engine and the Oracle database as its repository.

See Also: *Oracle Text Application Developer's Guide* and *Oracle9i Database Performance Tuning Guide and Reference*

Sufficient RAM Along with the resource requirements for the database and the Text indexing engine, also consider the memory requirements of the Ultra Search crawler. The Ultra Search crawler is a pure Java program. Out of the box, when the crawler is launched, the JVM is configured to start with 25MB and grow to 256MB. When crawling very large amounts of data, these values might need to be adjusted.

The Ultra Search administration tool is a J2EE 1.2 standard Web application. Therefore, it can be installed and run on a separate host from the Ultra Search backend (server component). However, you might want to install and run this on the same host as the Ultra Search backend. Regardless of your choice, allocate enough memory for the J2EE engine. Oracle recommends using the Oracle HTTP Server with the Oracle J2EE container. Allocate enough memory for the HTTP Server as well as the JDK that runs the J2EE engine.

Sufficient Disk Space Because customer requirements vary widely, Oracle cannot recommend a specific amount of disk space. However, as a general guideline, the minimal requirements are as follows:

- Approximately 3GB of disk space for the OracleAS infrastructure or database and the Ultra Search backend (server component).
- 15MB of disk space for the Ultra Search middle tier on top of the Web server's disk requirements.
- For each remote crawler host, the same amount of disk space as needed to install the Ultra Search backend.
- Disk space for a large TEMPORARY tablespace. As a general guideline, create a TEMPORARY tablespace as large as possible, depending on the RAM on your host.
- Disk space for the Ultra Search instance user's tablespace.
 - The Ultra Search instance user is a database user that you must explicitly create. All data that is collected and processed as part of crawling and indexing is stored in this user's schema.

- As a general guideline, create the tablespace as large as the total amount of data that you want to index. For example, if you approximate that the total amount of data to be crawled and indexed is 10GB, then create a tablespace that is at least 10GB for the Ultra Search instance user. Make sure to assign that tablespace as the default tablespace of the Ultra Search instance user.

Software Requirements

The Ultra Search middle tier components are Web applications. Therefore, they require a Web server to run. Oracle recommends using OracleAS.

Installing the Ultra Search Backend

The Ultra Search backend (server component) is included as part of the Oracle database. It is installed during the Oracle database installation. It is installed in the same Oracle home directory as the database server tier.

The Ultra Search backend consists of the following:

- Ultra Search data dictionary and PL/SQL packages
- Ultra Search crawler Java classes
- Ultra Search remote crawler
- Ultra Search product libraries

Installing and Configuring the Ultra Search Backend

Step 1: Install OracleAS Infrastructure and the Ultra Search Backend

Database Release Start up the Oracle Universal Installer (OUI) on the relevant host. After choosing the destination Oracle home name and full path, choose the option "Oracle9i Server." Ultra Search backend (server component) is installed with the Oracle database by default.

OracleAS Infrastructure Start up the OUI on the relevant host. After choosing the destination Oracle home name and full path, choose the option "OracleAS Infrastructure 10g" Ultra Search backend (server component) is installed with the OracleAS infrastructure by default.

During the installation of OracleAS or the Oracle database server, the Ultra Search backend is installed. The following activity occurs during this process:

- All Ultra Search backend files are copied into a directory named "ultrasearch". This directory resides immediately under the `ORACLE_HOME` of the designated database installation.
- The database user `WKSYS` with a randomized password is created. You should change this password later for security purposes. All Ultra Search database objects are installed in this user's schema.
- Various PL/SQL scripts are run against the database as user `WKSYS`. These scripts install and create various database objects.

See Also: "Changing Ultra Search Schema Passwords" on page 3-2 for information on changing the `WKSYS` password

See your installation guide for information on setting necessary environment variables.

Step 2: Set the Environment to Use the INSO Filter

The Ultra Search crawler uses the Oracle Text INSO filter `ctxhx`, which requires that your shared library path environment variable contain the `$ORACLE_HOME/ctx/lib` path. Without that, filtering fails for any binary document.

At installation, the Oracle Installer automatically sets the variable to include `$ORACLE_HOME/ctx/lib`. However, if, after the installation, you restart the database, then you must manually set your shared library path environment variable to include `$ORACLE_HOME/ctx/lib` before starting the Oracle process. You must restart the database to pick up the new value for filtering to work.

For example, on UNIX set the `$LD_LIBRARY_PATH` environment variable to include `$ORACLE_HOME/ctx/lib`, and on Windows set the `$PATH` environment variable to include `$ORACLE_HOME/bin`.

Step 3: Configure the Oracle Database for Ultra Search

After you have installed all Ultra Search components, you can optionally configure the Oracle database. This is a post-installation operation.

See Also: "Configuring the Oracle Server for Ultra Search" on page 3-2

Configure a Secure Ultra Search Installation

Step 1: Check the database version requirements and configure Oracle Identity Management.

Before you can set up a secure Ultra Search installation, you must do the following:

- Install or upgrade the Oracle database to 9.2.0.4 or higher. The middle tier and IM (identity management) version should be 9.0.4 or higher. You can use RepCA to convert a 9.2.0.4 database to an iAS 9.0.4 metadata repository.
- Install and configure the Oracle Internet Directory (OID)
- Configure the Oracle-OID link

Secure search functionality requires that the Ultra Search database is Oracle version 10.1.0 or higher and that the Ultra Search database is linked to a compatible instance of OID. This is necessary because Ultra Search utilizes XML DB functionality, which requires a certain version of Oracle. XML DB, in turn, requires a live link to OID, through which it retrieves all LDAP principal information. The Oracle-OID link must be running at all times for secure search to work. To set up this link, configure the Oracle Database to use Oracle Identity Management.

See Also: *Oracle9i Database Administrator's Guide* for details on configuring the database to use Oracle Identity Management and OID

Step 2: Restart the Oracle listener.

In the previous step, you configured the Oracle Database to use Identity Management. That process involved configuring the Oracle Home for directory usage. You must make sure to restart the Oracle listener to inherit the changes made to the Oracle Home. Restart the listener, if you have not already done so.

Step 3: Install or upgrade Ultra Search, if necessary.

After you have configured the Ultra Search database to work with OID, you can install or upgrade the Ultra Search backend (server component) into the Oracle Server, if you have not already done so.

Step 4: Create the `/sys/apps/ultrasearch` folder.

Immediately after installation or upgrade, you must run a SQL script to create the `/sys/apps/ultrasearch` folder in the XML DB repository. This folder stores all Ultra Search ACLs in XML DB.

To create the `/sys/apps/ultrasearch` folder, do the following:

1. cd to \$ORACLE_HOME/ultrasearch/admin
2. Login to the Ultra Search database using SQL*Plus as user WKSYS
3. Invoke the SQL script: @wk0prepxdb.sql

See Also: "Changing Ultra Search Schema Passwords" on page 3-2 for information on changing the WKSYS password

Upon termination, the wk0prepxdb.sql script lists all Ultra Search-related XML DB resources by running the following SQL:

```
SELECT any_path FROM resource_view WHERE any_path LIKE '%ultrasearch%';
```

Execution of that SQL statement must show two rows:

```
/sys/apps/ultrasearch  
/sys/apps/ultrasearch_acl.xml
```

If you do not see this confirmation, then this step has failed, and you cannot proceed. Recheck that all previous steps were performed correctly.

Step 5: Turn on secure search functionality in Ultra Search.

Because there is currently no way to programatically verify a proper Oracle-OID installation, the secure search functionality in Ultra Search is turned off by default. You must explicitly turn on this feature after completing all previous steps.

Step 6: Turn On Secure Search in the Query Application.

To turn on secure search functionality in Ultra Search:

1. Login to the Ultra Search database using SQL*Plus as user WKSYS
2. Invoke the following PL/SQL API: `exec WK_ADM.SET_SECURE_MODE(1)`

The argument (1) indicates that you are turning on secure search.

After you have turned on secure search functionality, you can create Ultra Search instances that are secure search-enabled.

Note: At any subsequent point in time, you can turn off security by invoking `WK_ADM.SET_SECURE_MODE(0)`. Doing so designates that any instances created after that will not support secure searches. However, existing secure search-enabled instances are not modified. Hence, if the Oracle-OID link ceases to function, you cannot perform searches on crawled documents that are secured.

Ultra Search supports secure searches, which return only documents satisfying the search criteria that the search user is allowed to view.

To turn on secure search in the query application, follow these steps:

1. Deploy Ultra Search query (`sample.ear`).
2. Edit the OC4J `jazn.xml` file to connect to OID. For example:

```
<jazn provider="LDAP" default-realm="us" location="ldap://localhost:3060">
<property name="ldap.user" value="orcladmin"/>
<property name="ldap.password" value="!welcome"/>
</jazn>
```
3. Restart OC4J.
4. Edit `applications/ultrasearch_query/META-INF/orion-application.xml` to turn on JAZN LDAP.
5. Edit `applications/ultrasearch_query/query/WEB-INF/web.xml` to enable login functionality in `usearch.jsp`. For example:

```
<init-param>
<param-name>login enabled</param-name>
<param-value>true</param-value>
</init-param>
```
6. Enable `mod_ossso` in Apache.
7. Access `http://<hostname>:<port>/ultrasearch/query/usearch.jsp` to see the login function, and test secure search.

See Also: "Secure Search" on page 1-6

Installing the Backend on an Existing Database or Metadata Repository

The Ultra Search backend (server component) can be installed on top of an existing Oracle 9*i* or later database. This can be done in two ways:

- Install Oracle Portal on to a database with the Oracle Portal Configuration Assistant (OPCA). This will also install Ultra Search.
- Use repCA to create an Oracle 10g Application Metadata Repository. As part of this process, Ultra Search is installed. This is probably the easiest way, but it comes with an overhead of having all OracleAS component schemas also installed on the target database.

Note: This feature is supported on the OracleAS release only.

Because the repCA approach is discussed in the "Using an Existing Database for the OracleAS Metadata Repository" section of the *Oracle Application Server 10g Installation Guide*, this section covers only the database install with OPCA.

Database Requirements

The database has the following requirements:

- The Java development kit (JDK) must be release 1.2.207 or higher.
- The Oracle database must be Oracle9*i* release 1 (9.0.1.4) or higher.
- The initialization parameter file must have `JOB_QUEUE_PROCESSES` set to two or higher.
- Oracle Text must be installed in the database. Oracle Text provides the text indexing and search capabilities required to index and query data retrieved from your data sources, such Web sites or database tables.

Note: To run the crawler, Ultra Search requires a Java 1.2 compliant runtime environment (JDK) to be installed on the database computer.

Installing the Ultra Search Backend on an Existing Oracle9*i* Database with OPCA

Follow these steps to install the backend (server component) on an existing Oracle9*i* database:

1. Launch the Oracle Portal Configuration Assistant (OPCA) to install Portal into the customer database. OPCA also installs the Ultra Search schema and database objects. If OPCA detects that the Ultra Search system schema `WKSYS` already exists, and the Ultra Search release is not the latest release, then it asks you to choose from the following three options:
 - **Deinstall/Reinstall:** This drops the existing `WKSYS` user, creates a new `WKSYS` user, and reloads the Ultra Search packages. All data collected by Ultra Search is deleted.
 - **Migrate:** This asks you to follow the instructions on the Ultra Search migration document for manual migration.

See Also: "Upgrading Ultra Search" on page 3-9
 - **Abort:** This stops OPCA from loading Ultra Search PL/SQL database packages into the existing database.

Note: You should immediately change the `WKSYS` password to avoid security problems. For details, see "Changing Ultra Search Schema Passwords" on page 3-2.

OPCA configures OracleAS to work with Ultra Search. Configuration involves the following files: `mod_oc4j.conf`, `server.xml`, and `ultrasearch.properties`. For details, see "Installing the Ultra Search Middle Tier on Web Server Hosts" on page 2-11.

2. Transfer files and configure the database. Ultra Search requires that certain files exist on the database tier's file system. Because the database tier might exist in a physically unreachable place, and because neither the Oracle Universal Installer nor OPCA have access to the remote database tier's file system, you must manually copy the necessary files. In the `ORACLE_HOME/ultrasearch/setup/` directory, there are several files including `setupDB.jar`, `setupDB.bat`, `setupDB.csh`, and `setupDB.sh`. Copy these files to the `ORACLE_HOME` of your remote database tier.
3. Set the files to the correct location by running the following scripts for your operating system.

For Windows: Edit `setupDB.bat` using any text editor. After the line "`set ORACLE_HOME=`" put the directory name where you installed the Oracle

database. Save the change and run `setupDB.bat` by double clicking `setupDB.bat` or by running `setupDB` in the command prompt.

For UNIX: Run `setupDB.sh` under the Bourne shell prompt. The system sets the current directory as `ORACLE_HOME` and asks for the path to your `JDK_HOME`. If the default path to the `jar` executable or Java executable is not correct, then the system asks for the path to `JDK`. If you do not have `JDK` installed, then download the latest `JDK` and run `setupDB.sh` again.

Note: To run this script, the executable bit of `setupDB.sh` should be turned on.

Configuring the Default Ultra Search Instance

The Ultra Search installer creates a default out of the box Ultra Search instance based on the default Ultra Search test user. So, you can test Ultra Search functionality based on the default instance after installation.

The default instance name is `WK_INST`. It is created based on the database user `WK_TEST`. The default user password is `WK_TEST`.

For security purposes, `WK_TEST` is locked after the installation. The administrator should login to the database as `DBA` role, unlock the `WK_TEST` user account, and set the password to be `WK_TEST`. (The password expires after the installation.) If the password is changed to anything other than `WK_TEST`, then you must also update the cached schema password using administration tool **Edit Instance** page after you change the password in the database.

The default instance is also used by the Ultra Search sample query application. Make sure to update the `data-sources.xml` file.

See Also:

- "Schema Password" on page 7-11
- "Editing the data-sources.xml File" on page 2-22

Installing the Ultra Search Middle Tier on Web Server Hosts

The Ultra Search middle tier includes the following:

- Ultra Search administration tool

- Ultra Search Java query API
- Ultra Search sample query applications

For the Oracle Application Server release, the Ultra Search middle tier is part of the Application Server installation. You must choose the "OracleAS Portal and Wireless" option from the Oracle Universal Installer menu to install and configure the Ultra Search middle tier with the Application Server install.

For the database release, the Ultra Search middle tier is installed with the Ultra Search backend (server component) during the database server install. This is also part of the database client install. The Ultra Search middle tier is installed and configured with Oracle J2EE container (OC4J).

See Also: *Oracle Application Server 10g Administrator's Guide* for information on how to change the Infrastructure Services (for example, a different Oracle Internet Directory or Metadata Repository) used by an Ultra Search middle tier

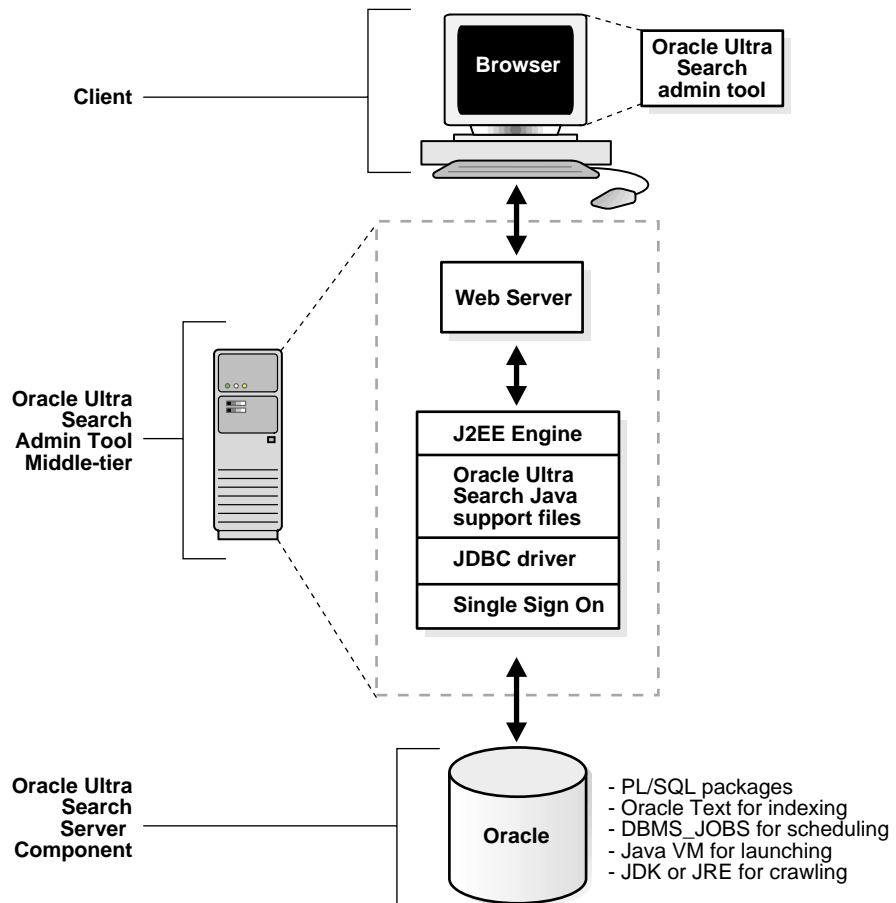
Web Applications Concepts

The Ultra Search administration tool and the Ultra Search query applications are J2EE-compliant Web applications. These are three tier architecture applications. Figure 2-1 shows the relationship between the browser (the first tier), the Web server and the servlet engine (the middle tier), and the Oracle Database (the third tier).

The Web server accepts requests from the browser and forwards the requests to the servlet engine for processing. The Ultra Search middle tier then communicates with the Oracle database through the JDBC, as in Figure 2-1.

You can use any browser to access the Ultra Search administration tool or Ultra Search sample query application. The URLs are described in the following section.

Figure 2-1 Ultra Search Architecture



Browser Requirements

To use the administration tool, your browser must be Netscape version 4.0 or Microsoft Internet Explorer version 4.0 or higher.

Installing the Middle Tier with the Oracle Database Release

Choose the installation option. Start up OUI on the relevant host. Choose to install the Oracle9*i* client. Make sure to choose the "Administrator Install" or the "Custom

Install" option. The OUI prompts for an Oracle home directory in which to install the middle tier. This directory is referred to as `$ORACLE_HOME`.

OUI automatically configures the Ultra Search middle tier with Oracle J2EE container (OC4J). Proceed to "Editing the ultrasearch.properties File" on page 2-24.

Installing the Middle Tier with the Oracle Application Server Release

Start the OUI on the relevant host. Choose the destination Oracle home name and full path, and complete the following steps:

1. Choose the option "OracleAS Application Server 10g", and click **Next**.
2. Choose the option "B. Portal and Wireless", and click **Next**.
3. On the "Configuration Options" screen, make sure "OracleAS Portal" is checked. This allows Oracle Portal Configuration Assistant (OPCA) to configure Oracle HTTP Server and OC4J with Ultra Search. If you uncheck this option, then you must follow the instructions under "Configuring the Middle Tier with Oracle HTTP Server and OC4J" to set up Oracle HTTP Server and OC4J manually.
4. Continue with the installation until OracleAS is successfully installed.

Note: If you decide to use a third party J2EE container or a servlet engine, then uncheck the option "OracleAS Portal" on the "Configuration Options" screen of Oracle Installer, and see the "Deploying the Ultra Search EAR File on a Third Party Middle Tier" on page 2-19. Upon completion of this step, all middle tier files are copied under the `$ORACLE_HOME`.

If you checked the "OracleAS Portal" option on the "Configuration Options" Oracle Installer screen, then the configuration steps in the following section are automatically performed by the Oracle Portal Configuration Assistant (OPCA). Proceed directly "Editing the data-sources.xml File" on page 2-22.

If not, then you must manually perform the steps under "Configuring the Middle Tier with Oracle HTTP Server and OC4J" on page 2-15 to configure your existing Web server.

You can also deploy Ultra Search Web applications using Oracle Enterprise Manager.

See Also:

- *Oracle9i Database Administrator's Guide* for more information on Enterprise Manager
- The Troubleshooting appendix in *Oracle Application Server 10g Installation Guide* for more information on OracleAS Configuration Assistants

Configuring the Middle Tier with Oracle HTTP Server and OC4J

Note: For Oracle Database, Oracle Containers for J2EE (OC4J) is configured by default. You can still configure the HTTP Server and OC4J, but they will be in a different Oracle home.

To deploy Ultra Search Web applications, you must have a J2EE 1.2 container. Oracle recommends using Apache Web server and OC4J.

See Also: "Deploying the Ultra Search EAR File on a Third Party Middle Tier" on page 2-19 if you use a third party J2EE container or servlet engine

1. For OC4J configuration, modify the following OC4J configuration files: `server.xml`, `application.xml`, and `default-web-site.xml` in `$ORACLE_HOME/j2ee/OC4J_Portal/config/`. The configuration of OC4J works with Ultra Search J2EE applications.

See Also: OracleAS Containers for J2EE documentation for more information on deploying EAR and WAR applications and for the more advanced functionality of OC4J

- For `server.xml`, under `<application-server>` tag, add the following:

```
<application name="UltrasearchAdmin" path="$ORACLE_
HOME/ultrasearch/webapp/ultrasearch_admin.ear" />
```

```
<application name="UltrasearchQuery" path="$ORACLE_
HOME/ultrasearch/sample.ear" />
```

```
<application name="UltrasearchPortlet" path="$ORACLE_
HOME/ultrasearch/webapp/ultrasearch_portlet.ear" />
```

Note: These lines let OC4J know that it must deploy the Ultra Search EAR file, as well as define where this EAR files is. `Ultrasearch_admin.ear` contains the Ultra Search administration tool Web application. The `sample.ear` file contains the sample query JSP pages. After OC4J deploys `sample.ear`, you can see the `$ORACLE_HOME/ultrasearch/sample` directory. Use the JSPs in this directory to create your own query Web pages. For more information on this directory, see "Testing the Ultra Search Sample Query Applications" on page 2-26.

- For `application.xml`, under `<orion-application>` tag, add the following:

```
<library path="$ORACLE_HOME/ultrasearch/lib/ultrasearch_query.jar" />
<library path="$ORACLE_HOME/ultrasearch/webapp/config" />
<library path="$ORACLE_HOME/jlib/uix2.jar" />
<library path="$ORACLE_HOME/jlib/share.jar" />
<library path="$ORACLE_HOME/jlib/regexp.jar" />
<library path="$ORACLE_HOME/lib/mail.jar" />
<library path="$ORACLE_HOME/lib/activation.jar" />
<library path="$ORACLE_HOME/lib/xmlparserv2.jar" />
<library path="$ORACLE_HOME/jdbc/lib/nls_charset12.zip" />
<library path="$ORACLE_HOME/jdbc/lib/classes12.jar" />
```

The preceding libraries are required for the Ultra Search administration tool and query Web applications to run.

Note: `$ORACLE_HOME/ultrasearch/webapp/config` contains the `ultrasearch.properties` file. For more information, see "Editing the `ultrasearch.properties` File" on page 2-24.

- For `default-web-site.xml`

Under `<web-site>` tag, add the following:

```
<web-app application="UltrasearchAdmin" name="admin"
root="/ultrasearch/admin" />

<web-app application="UltrasearchQuery" name="query"
root="/ultrasearch/query" />
```

```
<web-app application="UltrasearchPortlet" name="query"
root="/provider/ultrasearch" />
```

The preceding lines describe which Web application (WAR file) in the Ultra Search EAR files is deployed.

- The application field describes the application name. It should match the application name in `server.xml`.
- The name field describes the Web application name. This should match the WAR file name within the EAR file corresponding to the application.
- For root, specify the virtual path for this Web application. The virtual path is the path under the URL. For the administrative Web application, access it using
`http://hostname.domainname:port/ultrasearch/admin/`.

Note: The virtual path for a particular Web application is defined in three files: `default-web-site.xml`, `mod_oc4j.conf`, and `application.xml` in the `META-INF` directory of the EAR file. (The `META-INF` is created by extracting the EAR file.) You must modify the root attribute of `web-app` in `default-web-site.xml`, and the value enclosed by tag `context-root` in `application.xml` to change the virtual path point to each Web application.

2. Modify `modOC4J` configuration files. Add the following to `$ORACLE_HOME/Apache/Apache/conf/mod_oc4j.conf`:

```
Oc4jMount /ultrasearch/           OC4J_Portal
Oc4jMount /ultrasearch/*         OC4J_Portal
Oc4jMount /ultrasearch/query     OC4J_Portal
Oc4jMount /ultrasearch/query/*  OC4J_Portal
Oc4jMount /ultrasearch/ohw      OC4J_Portal
Oc4jMount /ultrasearch/ohw/*    OC4J_Portal
Oc4jMount /ultrasearch/admin_sso OC4J_Portal
Oc4jMount /ultrasearch/admin_sso/* OC4J_Portal
Oc4jMount /ultrasearch/admin    OC4J_Portal
Oc4jMount /ultrasearch/admin/*  OC4J_Portal
```

3. Ultra Search sample pages require JDBC connections to the database as the instance owner. Due to JServ limitations in the Oracle9i release, the user name, password and connection string used to create the JDBC connection are

hard-coded inside the sample JSP code. To configure the JSP to query a specific instance, edit the JSP source code, and replace the user name, password, and connection string values. All sample JSP source code is in the OC4J applications directory.

The following files contain user name, password, and connection string values:

- 9i/gsearch.jsp
- 9i/display.jsp
- 9i/gsearchf.jsp
- 9i/gutil.jsp
- 9i/mail.jsp

Note: The Oracle9i JSP files are being deprecated. It is not necessary to configure them if you do not plan to use them.

Configuring the Administration Tool with Single Sign-On Server

Note: Single sign-on is available only with the Oracle Identity Management infrastructure.

To configure the Ultra Search administration tool with the Oracle Single Sign-On (SSO) server, you must also follow these steps in addition to the configuration in "Configuring the Middle Tier with Oracle HTTP Server and OC4J" on page 2-15.

1. For OC4J configuration, modify the following OC4J configuration files:

application.xml and default-web-site.xml in \$ORACLE_HOME/j2ee/OC4J_Portal/config/.

- For application.xml, under <orion-application> tag, add the following:

```
<library path="$ORACLE_HOME/jlib/repository.jar" />
<library path="$ORACLE_HOME/jlib/jndi.jar" />
<library path="$ORACLE_HOME/jlib/ldapjclnt9.jar" />
<library path="$ORACLE_HOME/j2ee/home/jazn.jar" />
<library path="$ORACLE_HOME/j2ee/home/jaas.jar" />
```

- For default-web-site.xml, under <web-site> tag, add the following:

```
<web-app application="UltrasearchAdmin" name="admin"
```

```
root="/ultrasearch/admin_sso" />
```

2. Modify modOC4J configuration files. Add the following to mod_oc4j.conf:

```
Oc4jMount /ultrasearch/admin_sso/* OC4J_Portal
```

3. Confirm the following:

- \$ORACLE_HOME/Apache/Apache/conf/httpd.conf includes oracle_apache.conf
- \$ORACLE_HOME/Apache/Apache/conf/oracle_apache.conf includes ultrasearch.conf
- \$ORACLE_HOME/ultrasearch/webapp/config/ultrasearch.conf has the following content:

```
# add alias for ultra search online help and welcome page
Alias /ultrasearch/doc/ "/private/nli/ora9ias/ultrasearch/doc/"
Alias /ultrasearch/ "/private/nli/ora9ias/ultrasearch/sample/"

<IfModule mod_osso.c>
  <Location /ultrasearch/admin_sso>
    require valid-user
    authType Basic
  </Location>
</IfModule>
```

Deploying the Ultra Search EAR File on a Third Party Middle Tier

Because Ultra Search EAR files contain only Web applications (WAR files), they can be made to deploy on any J2EE 1.2 container. To do so, you must know the Ultra Search WAR file name, the predefined URL root, and the Java library required. The following section explains the Ultra Search EAR files that you deploy in a standard J2EE 1.2 container. It does not contain information on the configuration of each J2EE 1.2 container.

See Also:

- The documentation of the third party J2EE container for its configuration
- "Configuring the Middle Tier with Oracle HTTP Server and OC4J" on page 2-15

Deploying the Administration Tool The Ultra Search administration tool is a J2EE-compliant Web application (`$ORACLE_HOME/ultrasearch/webapp/ultrasearch_admin.ear`). You can use Enterprise Manager to deploy/undeploy this web application.

To see the file structure of `ultrasearch_admin`, run the following command:

```
jar -tvf ultrasearch_admin.ear
META-INF/
META-INF/application.xml
META-INF/orion-application.xml
admin.war
admin_sso.war
ohw.war
```

Deploying the Sample Query Applications Ultra Search sample query applications are Web applications contained in the `$ORACLE_HOME/ultrasearch/sample.ear` file. This file is already compliant to the J2EE 1.2 standard. You should not have to change this file to deploy it.

The following is the file structure of `sample.ear`. Extract the archived file by running the following command:

```
jar tf sample.ear

META-INF/application.xml
META-INF/orion-application.xml
query.war
welcome.war
rewriter/SampleRewriter.java
agent/SampleAgent.java
agent/README.html
```

All the query JSP pages are contained in `query.war`. This file is a servlet 2.2 compliant Web application. Deploy it alone with any servlet 2.2 engine. The context root for `query.war` is `/ultrasearch/query`. It is defined in the `META-INF/application.xml` of the `sample.ear` file. You can change it by editing this file.

The following are the Java libraries needed for Ultra Search sample query application:

```
$ORACLE_HOME/ultrasearch/webapp/config
$ORACLE_HOME/jdbc/lib/classes12.jar
$ORACLE_HOME/jdbc/lib/nls_charset12.zip
$ORACLE_HOME/ldap/jlib/ldapclnt9.jar
```

```
$ORACLE_HOME/lib/xmlparserv2.jar
$ORACLE_HOME/lib/activation.jar
$ORACLE_HOME/lib/mail.jar
```

Ultra Search query applications also use the connection pooling functionality of J2EE container. You must define a container authenticated data source. This data source must return an Oracle connection. Oracle recommends using the Java class equal to `oracle.jdbc.pool.OracleConnectionCacheImpl` for this data source.

In addition, the data source should contain the field location equal to `jdbc/UltraSearchPooledDS`, user name, password equal to the Ultra Search instance owner's database user name, and password and URL equal to the JDBC connection string in the form of `"jdbc:oracle:thin:@database_host:port:oracle_sid"`.

See Also: "Editing the data-sources.xml File" on page 2-22 for the data source configuration of the Oracle J2EE container

Deploying the Ultra Search Portlet Ultra Search Portlet is a Web application contained in the `$ORACLE_HOME/ultrasearch/webapp/ultrasearch_portlet.ear` file. This file is compliant to the J2EE 1.2 standard. This file is similar to `sample.ear` in terms of file structure. Extract the archived file by running the following command:

```
jar -xvf ultrasearch_portlet.ear

ultrasearch_portlet.ear
    META-INF/
        application.xml
    query.war
    agent/
    index.html
```

All the query JSP pages are contained in `query.war`. This file is a servlet 2.2 compliant Web application. You can deploy it alone with any servlet 2.2 engine. The context root for `query.war` is `/provider/ultrasearch/`. It is defined in the `META-INF/application.xml` of the `ultrasearch_portlet.ear` file. You can change it by editing this file.

The following Java libraries are needed for the Ultra Search Portlet:

```
$ORACLE_HOME/jdbc/lib/classes12.jar
$ORACLE_HOME/jdbc/lib/nls_charset12.zip
$ORACLE_HOME/lib/xmlparserv2.jar
$ORACLE_HOME/lib/activation.jar
$ORACLE_HOME/lib/mail.jar
```

Ultra Search Portlet uses the connection pooling functionality of J2EE container. You must define a container authenticated data source. This data source must return an Oracle connection. Oracle recommends using the Java class equal to `oracle.jdbc.pool.OracleConnectionCacheImpl` for this data source.

In addition, the data source should contain the field location equal to `jdbc/UltraSearchPooledDS`, user name, password equal to the Ultra Search instance owner's database user name, and password and URL equal to the JDBC connection string in the form of `jdbc:oracle:thin:@database_host:oracle_port:oracle_sid`.

See Also: "Editing the data-sources.xml File" on page 2-22 for the data source configuration of Oracle J2EE container

Editing the data-sources.xml File

Caution: Storing clear text passwords in `data-sources.xml` poses a security risk. Avoid this by using password indirection to specify the password. This lets you enter the password in `jazn-data.xml`, which automatically gets encrypted, and point to it from `data-sources.xml`. For more information, see "Creating An Indirect Password" in *Oracle Application Server Containers for J2EE Security Guide*.

Oracle Application Server Infrastructure

The Ultra Search Oracle Application Server query API uses the data source functionality of the J2EE container. Under directory `$ORACLE_HOME/j2ee/OC4J_Portal/config`, edit the file `data-sources.xml`. Under tag `<data-sources>` add the following:

```
<data-source
  class="oracle.jdbc.pool.OracleConnectionCacheImpl"
  name="UltraSearchDS"
  location="jdbc/UltraSearchPooledDS"
  username="username"
  password="password"
  url="jdbc:oracle:thin:@database_host:oracle_port:oracle_sid"
/>
```


Where *username* and *password* are the Ultra Search instance owner's database user name and password, *database_host* is the host name of the back end database computer, *oracle_port* is the port to the user's Oracle database, and *oracle_sid* is the SID of the user's Oracle database. In addition to user name, password, and JDBC URL, `data-sources.xml` also allows configuration of the connection cache size, as well as the cache scheme.

The following tag specifies the minimum and maximum limits of the cache size, the inactivity time out interval, and the cache scheme.

If you are adding the data source for the default Ultra Search instance user `wk_test`, then make sure to unlock `wk_test` first.

See Also: "Configuring the Default Ultra Search Instance" on page 2-11

```
<data-source
  class="oracle.jdbc.pool.OracleConnectionCacheImpl"
  name="UltraSearchDS"
  location="jdbc/UltraSearchPooledDS"
  username="wk_test"
  password="wk_test"
  url="jdbc:oracle:thin:@localhost:5521:isearch"
  min-connections="3"
  max-connections="30"
  inactivity-timeout="30">
  <property name="cacheScheme" value="1"/>
</data-source>
```

Note: The URL of the JDBC data source can be provided in the form of `jdbc:oracle:thin:@[hostname]:[port]:[sid]` or in the form of a TNS keyword-value syntax, such as

```
"jdbc:oracle:thin:@(DESCRIPTION=(LOAD_
BALANCE=yes)(ADDRESS_LIST=(ADDRESS=(PROTOCOL=TCP)
(HOST=cls02a)(PORT=3999))(ADDRESS=(PROTOCOL=TCP)
(HOST=cls02b)(PORT=3999)))(CONNECT_DATA=(SERVICE_
NAME=acme.us.com)))"
```

There are three types of caching schemes:

- `DYNAMIC_SCHEME = 1`
- `FIXED_WAIT_SCHEME = 2`

- `FIXED_RETURN_NULL_SCHEME = 3`

See Also: *Oracle Application Server Containers for J2EE Security Guide*

Database Release

For the database release, follow the same directions as earlier. However, the `data-sources.xml` file is slightly different:

```
<data-source
class="oracle.jdbc.pool.OracleDataSource"
name="UltraSearchDS"
location="jdbc/UltraSearchPooledDS"
username="wk_test"
password="wk_test"
url="jdbc:oracle:thin:@dlsun1517.us.oracle.com:5521:dczade4"
connectionCachingEnabled="true"/>
```

Editing the `ultrasearch.properties` File

The `$ORACLE_HOME/ultrasearch/webapp/config/ultrasearch.properties` file contains configuration information used by Ultra Search middle tier. You do not need to edit this file, because it is automatically configured by the Oracle installer.

Here is an example of the `ultrasearch.properties` file:

```
connection.driver=oracle.jdbc.driver.OracleDriver
connection.url=jdbc:oracle:thin:@ldap://dlsun8888.cn.oracle.com:3060/iasdb,cn=oraclecontext
oracle.net.encryption_client=REQUESTED
oracle.net.encryption_types_client=(RC4_56,DES56C,RC4_40,DES40C)
oracle.net.crypto_checksum_client=REQUESTED
oracle.net.crypto_checksum_types_client=(MD5)
oid.app_entity_cn=m16bi.sgtcnsun03.cn.oracle.com
domain=us.oracle.com
```

Where

- `connection.driver` specifies the JDBC driver you are using.
- `connection.url` specifies the database to which the middle tier connects. Ultra Search supports following formats:

- *host:port:SID* (where *host* is the full host name of the Oracle base instance running Ultra Search, *port* is the listener port number for the Oracle Database instance, and *SID* is the Oracle Database instance ID)
- HA-aware string (for example, TNS keyword-value syntax)

Here is an example `connection.url` string:

```
connection.url=jdbc:oracle.thin:@ultrasearch.us.oracle.com:1521:myInstance
```

- `oracle.net.encrypted_client`, `oracle.net.encrypted_types_client`, `oracle.net.crypto_checksum_client`, and `oracle.net.crypto_checksum_types_client` control the properties of the secure JDBC connection made to the database. See *Oracle9i JDBC Developer's Guide and Reference* for more information.
- `oid.app_entity_cn` specifies the Ultra Search middle tier application entity name.
- `domain` specifies the common domain for the IM (identity management) machine and the Ultra Search middle tier machine. This enables delegated administrative service (DAS) list of values to work with Internet Explorer. For example, if the Ultra Search middle tier is in `us.oracle.com` and the IM machine is `uk.oracle.com`, then the common domain is `oracle.com`. Add the following line in `ultrasearch.properties`: `domain=oracle.com`

Starting the Web Server

With the OracleAS release, start the Web server using the Oracle Enterprise Manager Application Server Control.

See Also: *Oracle Application Server 10g Administrator's Guide* for information on the Application Server Control

With the database release, do the following:

```
java -jar $ORACLE_HOME/j2ee/home/oc4j.jar -config
$ORACLE_HOME/j2ee/OC4J_SEARCH/config/server.xml
```

Testing the Ultra Search Administration Tool

Check that the Web Server is running.

Test your changes by attempting to log on to the administration tool:

- **Visit:**
`http://hostname.domainname:port/ultrasearch/admin/index.jsp`
where `hostname.domainname` is the full name of the host where you have installed the Ultra Search middle tier, and `port` is the default Web server port.
- During the installation of the Ultra Search backend (server component), you should have created a new Ultra Search instance owner. Log on to the Ultra Search administration tool by entering the Ultra Search instance owner's database user name and password.
- The nature of JSP pages is such that the first time any page is accessed, it takes a few seconds to compile. Subsequent accesses are much faster.
- If you log on to the Ultra Search administration tool successfully, then you have completed the Ultra Search administration tool configuration process.

Testing the Ultra Search Sample Query Applications

After you verify that the Ultra Search administration tool is working, you should be able to run the Ultra Search sample query applications.

To test the Ultra Search sample query applications, do one of the following:

- **Visit**
`http://hostname.domainname:port/ultrasearch/query/search.jsp`
- **Follow the links in the Ultra Search welcome page:**
`http://hostname.domainname:port/ultrasearch/index.html`

See Also: "Configuring the Middle Tier with Oracle HTTP Server and OC4J" on page 2-15 for information about configuring the JSP to query a specific instance

Locations for sample query applications are listed in the following section. Access the sample query source code by going to the directories list. You can also see a working demo of each sample query JSP page with the URL root, and you can append the correct JSP file name at the end of the URL root.

The root query directory is `$ORACLE_HOME/ultrasearch/sample/query/`.

The URL root for the query is

`http://hostname.domainname:port/ultrasearch/query/`.

The OracleAS query (query sample pages that use the OracleAS query API and include `usearch.jsp` and `search.jsp`) is in `$ORACLE_HOME/ultrasearch/sample/query/`.

The URL root for the OracleAS query is in

`http://hostname.domainname:port/ultrasearch/query/.`

(For example: `access search.jsp` with

`http://hostname.domainname:port/ultrasearch/query/search.jsp.`)

The Oracle Database query (query JSP that uses the 9i query API and includes `gsearch.jsp`) is in `$ORACLE_HOME/ultrasearch/sample/query/9i/.`

The URL root for the Oracle Database query is in

`http://hostname.domainname:port/ultrasearch/query/9i/.`

Portlet is in `$ORACLE_HOME/ultrasearch/sample/query/portlet/.`

The URL root for Portlet is in

`http://hostname.domainname:port/ultrasearch/query/portlet/.`

Taglib is in `$ORACLE_HOME/ultrasearch/sample/query/tag/.`

The URL root for taglib is in

`http://hostname.domainname:port/ultrasearch/query/tag/.`

Installing the Backend on Remote Crawler Hosts

The Ultra Search remote crawler allows multiple crawlers to run in parallel on different hosts. However, all remote crawler hosts must share common resources, such as common directories and a common Ultra Search database.

Installing the Backend on Remote Crawler Hosts

The Ultra Search remote crawler is part of the Ultra Search backend (server component). Therefore, the installation procedure is the same as installing the Ultra Search backend

On each remote crawler host, the Ultra Search backend is installed under a common directory known as the Oracle home. You should have been prompted by the Oracle Universal Installer to enter this directory. The Oracle home directory is referred to as `$REMOTE_ORACLE_HOME`.

If you choose not to install the Oracle HTTP Server during the OracleAS installation, then you must perform the following steps manually for remote crawling:

- Locate `$REMOTE_ORACLE_HOME/ultrasearch/tools/remotecrawler/scripts/unix/define_env` on a UNIX system or `$REMOTE_ORACLE_`

HOME/ultrasearch/tools/remotecrawler/scripts/winnt/define_env.bat on a Windows system.

- Replace %ORACLE_HOME% with the value of the REMOTE_ORACLE_HOME environment variable.
- Replace %s_jreLocation% with the directory path of a Java runtime environment (JRE) version 1.2.2 and higher. You should specify the root directory of the JRE.
- Replace %s_jreJDBCclassfile% with the full path and file name of the Oracle JDBC Thin driver (version 12).

Configuring the Backend on Remote Crawler Hosts

The only configuration needed for an Ultra Search remote crawler host is to register the host with the Ultra Search system. The registration process is done by running a SQL script on the Ultra Search remote crawler host. The SQL script connects over SQL*Plus to the OracleAS middle tier and registers the remote crawler host.

1. Locate the correct Oracle home.

The Ultra Search middle tier is installed under a common directory known as the Oracle home. If you have installed other Oracle products prior to the Ultra Search middle tier, then you could have multiple Oracle homes on your host. The registration script requires that you enter the Oracle home directory in which the Ultra Search middle tier is installed.

2. Locate the WKSYS super-user password.

You must run the registration script as the WKSYS super-user or as a database user that has been granted super-user privileges.

3. Start SQL*Plus.

The SQL script is located in

```
/ultrasearch/tools/remotecrawler/scripts/common/register.sql  
1 under $REMOTE_ORACLE_HOME.
```

Be sure to run the correct version of SQL*Plus, because multiple versions can reside on the same host if you have previously installed some Oracle products. On UNIX platforms, make sure that the correct values for PATH, ORACLE_HOME and TNS_ADMIN variables are set. On Windows platforms, choose the correct menu item from the Start menu.

After you have identified how to run the correct SQL*Plus client, you must log on to the Ultra Search database. To do this, you might need to configure an Oracle Net service setting for the Ultra Search database.

See Also: *Oracle9i Net Services Administrator's Guide* for information on how to configure a service setting

After SQL*Plus is running, log on to the database using the schema and password that you located in Step 2.

4. Invoke the registration script.

Starting up SQL*Plus as the WKSYS super-user and enter the following:

```
@full_path_of_registration_script
```

For example, if value for \$REMOTE_ORACLE_HOME on a UNIX host is /home/oracle9i, then enter the following at the SQL*Plus prompt:

```
@/home/oracle9i/ultrasearch/tools/remotecrawler/scripts/unix/register.sql
```

If you are running SQL*Plus on Windows, and \$REMOTE_ORACLE_HOME is in d:\Oracle\Oracle9i, then enter the following at the SQL*Plus prompt:

```
@d:\Oracle\Oracle9i\ultrasearch\tools\remotecrawler\scripts\winnt\register.sql
```

The registration script prompts you for two variables. The following is a list of the variables and their descriptions:

REMOTE_CRAWLER_HOSTNAME: The DNS host name of the remote crawler host.

ORACLE_HOME: The Oracle home located in Step 1. For example, /u01/oracle9i on a UNIX host or D:/u01/oracle9i on a Windows host. (Be careful to use forward slashes for Windows hosts.)

The registration script invokes the wk_crw.register_remote_crawler PL/SQL API. The REMOTE_CRAWLER_HOSTNAME and ORACLE_HOME variables are used to compose arguments for the wk_crw.register_remote_crawler API.

5. Verify and complete the remote crawler profile configuration.

Be sure to enter the correct values for both variables. To verify that the registration has completed correctly, log on to the Ultra Search administration tool. Click the **Remote Crawler Profiles** subtab in the **Crawler** tab. You should see the host name of the remote crawler host you have just registered in the

remote crawler profile list. Click **Edit** to complete the configuration process for the remote crawler profile.

Unregistering a Remote Crawler

If you enter any wrong values for the `register.sql` script, then you must unregister the remote crawler using the `unregister.sql` script. Invoke the `unregister` script the same way as you invoke the registration script. The `unregister.sql` script calls the `wk_crw.unregister_remote_crawler` PL/SQL API. After you have successfully unregistered the remote crawler, you can rerun the `register.sql` script.

Configuring Ultra Search in a Hosted Environment

Ultra Search is configured to be non-hosted during the default install. To change to a hosted environment, perform the following steps to configure Ultra Search in the hosted environment.

Preconfiguration Tasks for a Hosted Environment

Make sure the hosting mode is enabled. Also, make sure the subscriber is created in the OID server.

See Also: *OracleAS Portal Configuration Guide* section F.2 **Enabling Hosting on an Out-of-Box Portal** for instructions on how to enable the hosting mode, and section F.4 **Adding Subscribers** for instructions on how to add a subscriber to the SSO/OID

Configuring Ultra Search in the Subscriber Context

For each subscriber, run the following scripts to configure Ultra Search in the OID subscriber context. The script does the following:

- Creates the reference objects in the subscriber context.
- Creates default privilege group entry in the subscriber context.
- Updates the subscriber information in the Ultra Search metadata repository.

Script usage:

```
ORACLE_HOME/ultrasearch/setup/usca.sh -action add_subscriber -user <the OID user DN> -password <password of the user 'orcladmin'> -subscriber <the DN of the subscriber>
```


The OID user must have the 'iASAdmins' privilege. Before you run the script, make sure you have the execute permission on the script, and setup the ORACLE_HOME environment variable.

The following example configures Ultra Search in the subscriber 'dc=us, dc=oracle, dc=com':

```
ORACLE_HOME/ultrasearch/setup/usca.sh -action add_subscriber -user  
'cn=orcladmin' -password welcome1 -subscriber 'dc=us,dc=oracle,dc=com'
```

To drop the subscriber, first perform the following script to remove Ultra Search entries from the OID subscriber context:

```
ORACLE_HOME/ultrasearch/setup/usca.sh -action remove_subscriber -user <the OID  
user DN> -password <password of the user 'orcladmin'> -subscriber <the DN of the  
subscriber>
```

Post-Installation Information

This chapter contains the following topics:

- Changing Ultra Search Schema Passwords
- Configuring the Oracle Server for Ultra Search
- Managing Stoplists
- Upgrading Ultra Search
- Configuring the Query Application

Changing Ultra Search Schema Passwords

There are two Ultra Search system schemas created during installation: `WKSYS` and `WKPROXY`. You can update the schema password in the following way:

For the Oracle Database Release:

After the database is installed, all user schema accounts are locked. To log on as user `WKSYS` (or `WKPROXY`), unlock `WKSYS` (or `WKPROXY`) by running the following statement as the `SYSTEM` or `SYS` database user:

```
ALTER USER WKSYS ACCOUNT UNLOCK;
```

For the Oracle Application Server or the Oracle Collaboration Suite Release:

After the infrastructure database is installed, all user schema passwords are randomized. To log on as user `WKSYS` (or `WKPROXY`), change the `WKSYS` (or `WKPROXY`) schema password by following the link **Change Schema Password** from the Oracle Enterprise Manager **Infrastructure** page.

See Also: *Oracle Enterprise Manager Basic Installation and Configuration*

Configuring the Oracle Server for Ultra Search

The operations described in this section are database administration operations. They can be performed using Oracle Enterprise Manager or SQL*Plus.

Step 1: Tune the Oracle Database

Increase the Size of the Oracle Redo Logs, if necessary

Every instance of an Oracle database has an associated online redo log, which is a set of two or more online log files that record all committed changes made to the database. Online redo logs protect the database in the event of an instance failure. The size of redo log files determines the frequency of redo log file switches. This, in turn, significantly impacts text indexing speed. To reduce the frequency of log file switches, ensure that the redo log files are each 100MB or more.

The following section lists some tips on how to increase the redo log file sizes, if necessary. Enter the statements in the following section with the appropriate Oracle administrator privileges.

See Also:

- *Oracle9i Database Performance Tuning Guide and Reference*
- *Oracle9i Database Administrator's Guide*

1. Locate redo log files and determine their sizes:

```
SELECT v$logfile.member, v$logfile.group#, v$log.status, v$log.bytes
FROM v$log, v$logfile
WHERE v$log.group# = v$logfile.group#;
```

2. Add larger redo log files:

```
ALTER DATABASE ADD LOGFILE 'redo_log_directory/newredo1.log' size 100m;
ALTER DATABASE ADD LOGFILE 'redo_log_directory/newredo2.log' size 100m;
ALTER DATABASE ADD LOGFILE 'redo_log_directory/newredo3.log' size 100m;
```

A production database should have more log members for each log group, and different storage devices should be used to increase performance and reliability.

3. Drop the old log files. For each old redo log file, enter the ALTER SYSTEM SWITCH LOGFILE statement until that log file's status is INACTIVE. This is necessary to ensure that Oracle is not using that log file when you try to drop it.

Then, drop the old redo log file with the following statement:

```
ALTER DATABASE DROP LOGFILE 'redo_log_directory/redo01.log';
ALTER DATABASE DROP LOGFILE 'redo_log_directory/redo02.log';
ALTER DATABASE DROP LOGFILE 'redo_log_directory/redo03.log';
```

4. Manually delete the old log files from the file system. For each old redo log file, use the appropriate operating system statement to delete the unwanted log file from the file system.**Increase the Size of the Undo Space**

Every Oracle database must have a method of maintaining information that is used to roll back, or undo, changes to the database. Such information consists of records of the actions of transactions, primarily before they are committed. Oracle refers to these records collectively as undo. The undo space created by the Oracle Installer may be too small. Oracle recommends that you use automatic undo management and increase the undo space.

See Also: *Oracle9i Database Administrator's Guide* for details on using automatic undo management

Tune Oracle Initialization Parameters

Set the following values in the initialization file:

- `PROCESSES`: Set this to 50 or more.
- `SORT_AREA_SIZE`: Set this to 5MB or more.
- `SORT_AREA_RETAINED_SIZE`: Set this to 5MB or more.
- `JOB_QUEUE_PROCESSES`: Set this to three or higher. (Set it to at least one.) This is needed because the Ultra Search crawler is launched by scheduling a database job. If this is zero, then no database jobs are run. As a result, any attempts to launch the Ultra Search crawler fail. Also consider other requirements for job queue processes when you set this value.

For the latest information on initialization parameters relating to Ultra Search, see the Ultra Search Readme.

Step 2: Create and Assign the Temporary Tablespace to the CTXSYS User

The starter database created by the Oracle Installer may create a temporary tablespace that is too small. Oracle Ultra Search uses the Oracle Text engine intensively. Therefore, a large temporary tablespace must be created for the Oracle Text system user `CTXSYS`. If you want greater read and write performance, create the tablespace on raw devices.

When you have created the temporary tablespace, assign it as the temporary tablespace for the `CTXSYS` user. To do so, you must log on as the `SYSTEM` or `SYS` user. Assign the temporary tablespace to the `CTXSYS` user with the following statement:

```
ALTER USER CTXSYS TEMPORARY TABLESPACE new_temporary_tablespace;
```

See Also: *Oracle9i Database Administrator's Guide* for information on how to create a temporary tablespace

Step 3: Create a Large Tablespace for Each Ultra Search Instance User

For each Ultra Search instance, you must create a tablespace large enough to contain all data obtained during the crawling and indexing processes. This amount is subject to the amount of data you intend to crawl and index. However, it is often not possible to know in advance how much data you intend to collect. Try to obtain an estimate of the cumulative size of all data you want to crawl.

If you cannot estimate the size, then try to allocate as much space as possible. If you run out of disk space, then Ultra Search is able to resume crawling after you add more datafiles to the instance tablespace.

Here is an example of how to create a new tablespace:

```
CREATE TABLESPACE lmtbsb DATAFILE '/u02/oracle/data/lmtbsb01.dbf' SIZE 150M;
```

Pay attention to the `STORAGE` clause in your `CREATE TABLESPACE` statement. The amount of data to be stored in the tablespace can be very large. This can cause the Oracle server to progressively allocate many new extents when more storage space is needed. If the extent management clause specifies that each new extent is to be larger than the previous extent (that is, the `PCTINCREASE` setting is nonzero), then you could encounter the situation where the next extent that the Oracle server wants to allocate is larger than what is available. In such a situation, indexing halts until new extents can be added to the tablespace.

To mitigate this problem, certain instance-specific tables have explicit storage parameter settings. The initial extent size, next extent size, and `PCTINCREASE` setting are defined for these tables. These tables are created when a new instance is created. The tables and their storage clause settings are as follows:

```
DR$WK$DOC_PATH_IDX$I
      (initial extent size 5M, next extent size 50M, PCTINCEASE 1)
DR$WK$DOC_PATH_IDX$K
      (initial extent size 5M, next extent size 50M, PCTINCEASE 1)
```

If you want greater read and write performance, create the tablespace on raw devices.

Be sure to create a new large tablespace for each Ultra Search instance user.

See Also:

- *Oracle9i SQL Reference* for more information on creating tablespaces and managing storage settings
- *Oracle9i Database Administrator's Guide* for information on how to create a tablespace

Step 4: Create and Configure New Database Users for Each Ultra Search Instance

Ultra Search uses Oracle's fine grained access control feature to support multiple Ultra Search instances within one physical database. This is especially useful for large organizations or application service providers (ASPs) that want to host multiple disjoint search indexes within one installation of Oracle.

Note: Ultra Search requires that each Ultra Search virtual instance belong to a unique database user. Therefore, as part of the installation process, you must create one or more new database users to own all data for your Ultra Search instance.

If you intend to create more than one database instance, you should also create multiple user tablespaces: one for each user.

You must grant the `WKUSER` role to database users hosting new Ultra Search instances.

See Also: "Users Page" on page 7-42

Enter the following statements to create and configure a new user. Run these statements as the `WKSYS`, `SYSTEM`, or `SYS` database user.

```
CREATE USER username
        IDENTIFIED BY password DEFAULT TABLESPACE default_tbs
        TEMPORARY TABLESPACE temporary_tbs QUOTA UNLIMITED
        ON default_tbs;
```

where *username* = name of the Ultra Search instance owner

and *password* = password of the Ultra Search instance owner

and *default_tbs* = default tablespace for the Ultra Search instance created in step 3

and *temporary_tbs* = temporary tablespace created in step 2

```
GRANT WKUSER TO username;
```

After these steps are completed, `WKSYS` or an Ultra Search super-user can create an Ultra Search instance on this user schema.

If you want this user to have the general administrative privilege or the super-user privilege, then log on as an Ultra Search super-user or `WKSYS` and click go to the **Users** page in the administration tool to grant the appropriate privilege.

Step 5: Alter the Index Preferences

This step is optional.

An empty index is created when an Ultra Search instance is created. The existing index preferences, such as language-specific parameters, are defined in the `$ORACLE_HOME/ultrasearch/admin/wk0pref.sql` file.

You can modify these preferences so that all new Ultra Search instances use the modified preferences, or you can alter the index using your own preferences immediately after an instance is created. Alter the index using SQL.

Note: The crawler transforms all documents into HTML files with binary document filtering before indexing begins.

See Also:

- *Oracle Text Application Developer's Guide*
- *Oracle Text Reference*

Managing Stoplists

Every Ultra Search instance has a stoplist associated with it. A stoplist is a list of words that are ignored during the indexing process. These words are known as stopwords. Stopwords are not indexed because they are deemed not useful, or even disruptive, to the performance and accuracy of indexing.

Default Ultra Search Stoplist

During the installation process, a default stoplist is created for the Ultra Search product. Subsequently, when an Ultra Search instance is created, a copy of the default stoplist is created for the Ultra Search instance.

The default stoplist is created under the `WKSYS` schema. The default stoplist name is `wk_stoplist`. (This list is defined in the file `$ORACLE_HOME/ultrasearch/admin/wk0pref.sql`, which is run at installation).

Modifying Instance Stoplists

Modify the default stoplist by adding or removing stopwords from it. However, remember that these modifications do not affect existing Ultra Search instances. They only affect Ultra Search instances that are created after the modifications are made.

Modifying instance stoplists should be done as a last resort. Use one of the following methods:

- Modify the default stoplist before creating the instance.
- Replace the instance stoplist immediately after creating the instance.

Replacing the instance stoplist immediately after creating the instance affects only that instance. You must first create a user-defined stoplist.

In both cases, the result is that the Ultra Search instance stoplist is modified and defined before initial crawling. This means that all documents collected by the Ultra Search crawler are evaluated against the correct stoplist. It is important to modify the stoplist before initial crawling to avoid having to recrawl all documents again.

Modifying Instance Stoplists Before Initial Crawling

1. Modify the default stoplist before creating the instance:

For example, to add the stopword "web" to the default stoplist, log on as user WKSYS in SQL*Plus, and run the following statement:

```
EXEC ctx_ddl.add_stopword('wk_stoplist','web');
```

To remove the stopword "web" from the default stoplist, log on as user WKSYS in SQL*Plus, and run the following statement:

```
EXEC ctx_ddl.remove_stopword('wk_stoplist','web');
```

Subsequently, the stoplists of all new instances reflect the modifications made to the default stoplist.

2. Replace the instance stoplist immediately after creating the instance:

You must create a new user-defined stoplist. Log on as the owner of the instance in SQL*Plus, and run the following statements:

```
BEGIN  ctx_ddl.create_stoplist('example_stoplist');
        ctx_ddl.add_stopword('example_stoplist','example_stopword');
        ... (add more stopwords by repeated the previous
            line with new stopwords) ...
END;
/
```

To replace an instance stoplist with this new stoplist, log on as the owner of the instance in SQL*Plus, and run the following statement:

```
ALTER INDEX wk$doc_path_idx rebuild parameters('replace stoplist example_
```

```
stoplist');
```

See Also: "Changing Ultra Search Schema Passwords" on page 3-2 for information about changing the WKSYS password

Modifying Instance Stoplists After Initial Crawling

If necessary, alter an instance stoplist after initial crawling with one of the following methods:

1. Add stopwords to the instance stoplist:

Choosing to add stopwords to the instance stoplist does not affect any documents already crawled or indexed. This operation is not an expensive operation.

For example, to add the stopword "web" to the instance stoplist, log on as the owner of the instance in SQL*Plus, and run the following statement:

```
ALTER INDEX wk$doc_path_idx rebuild parameters('add stopword web');
```

2. Replace the instance stoplist after initial crawling:

Defining a new stoplist and replacing the instance stoplist with it invalidates the entire index. If you choose this method, you must force the Ultra Search crawler to recrawl all documents in the index. To do this, click **Process All Documents** in the **Edit Schedule** page. This is a very expensive operation. Therefore, this option should be the last resort.

Upgrading Ultra Search

Ultra Search is shipped with the Oracle Database, the Oracle Application Server, and the Oracle Collaboration Suite. To upgrade Ultra Search from a previous release to the most recent release, you must apply different procedures based on the product you are using.

This section contains the following topics:

- Pre-Upgrade Steps
- Upgrading Ultra Search Shipped with Oracle Database
- Upgrading Ultra Search Shipped with Oracle Application Server
- Upgrading Ultra Search Shipped with Oracle Collaboration Suite
- Upgrading Ultra Search to Oracle Collaboration Suite Release 1

See Also: "Ultra Search Release Information" describes the Ultra Search release numbering

Pre-Upgrade Steps

Before you upgrade, log on to the Ultra Search administration tool. Stop and disable all crawler synchronization schedules in every Ultra Search instance. You can enable all crawler synchronization schedules after the upgrade. See "Schedules Page" on page 7-32 for details on how to stop and disable the synchronization schedule.

Upgrading Ultra Search Shipped with Oracle Database

To upgrade Ultra Search shipped with the Oracle Database release, do the following:

1. Run the Ultra Search backend (server component) upgrade. This includes upgrading the Ultra Search database schemas and server files. Install the new Oracle software, and run Oracle Database Upgrade Assistant to upgrade the database and Ultra Search component to the new release. See the *Oracle Database Upgrade Guide* for details.
2. Follow the steps in "Installing the Ultra Search Middle Tier on Web Server Hosts" on page 2-11 to install the new Ultra Search middle tier.

Upgrading Ultra Search Shipped with Oracle Application Server

To upgrade Ultra Search shipped with the Oracle Application Server, do the following:

1. Install the new Oracle Application Server and use Oracle Application Server Upgrade Assistant to upgrade the middle tier. See the *Oracle Application Server 10g Upgrading to 10g (9.0.4)* section "Upgrading the Middle Tier" for details.
2. Perform the upgrade on the Ultra Search schema in the Oracle Application Server Metadata Repository. See the *Oracle Application Server 10g Upgrading to 10g (9.0.4)* section "Upgrading the Metadata Repository->Executing the Oracle Ultra Search Schema Upgrade Script" for details.

Upgrading Ultra Search Shipped with Oracle Collaboration Suite

If you are using the Oracle Collaboration Suite release 1 and want to upgrade to the most recent Oracle Collaboration release, then install the latest Oracle Collaboration Suite release and use Oracle Collaboration Suite Upgrade Assistant to upgrade both

the Ultra Search middle tier and backend (server component). See the “Oracle Collaboration Suite Installation and Configuration Guide” for details.

If you are using Ultra Search 9.0.2 (shipped with Oracle Application Server release) or Ultra Search 1.0.3 or 9.2 (shipped with Oracle Database release) and want to upgrade to the most recent Oracle Collaboration release, then perform the following upgrade procedures:

1. Get the Oracle Collaboration Suite release 1 software and upgrade your Ultra Search to Oracle Collaboration Suite release 1 first. See section "Upgrading Ultra Search to Oracle Collaboration Suite Release 1" on page 3-11 for details.
2. Install the latest Oracle Collaboration Suite release and use Oracle Collaboration Suite Upgrade Assistant to upgrade both Ultra Search middle tier and backend. See *Oracle Collaboration Suite Installation and Configuration Guide* for details.

Upgrading Ultra Search to Oracle Collaboration Suite Release 1

Ultra Search supports the following upgrades:

- Upgrade from Ultra Search 1.0.3 to 9.0.3
- Upgrade from Ultra Search 9.0.2 to 9.0.3
- Upgrade from Ultra Search 9.2 to 9.0.3

Upgrade is based on the backend (server component) only. Upgrade on the middle tier is not supported. Install the 9.0.3 middle tier in a separate Oracle home.

Upgrade from Ultra Search 1.0.3 to 9.0.3

Upgrading from Ultra Search 1.0.3 (Oracle9i Database 9.0.1) to 9.0.3 requires running the upgrade script and performing some manual steps.

The Ultra Search upgrade script first verifies the version of the current system, then upgrades the system and migrates user data. User data includes all dictionary and table data, such as information about the metadata, data sources, mappings, crawler schedules, authentication, and query statistics.

All crawler schedules and jobs created in the older version are disabled before data and system migration. **When migration is complete, the system administrator should re-activate the crawling schedule to re-index the document.** You do not need to reconfigure the system or re-enter any data. You can still query documents that were crawled and indexed by the previous version.

See Also: "Installing the Backend on an Existing Database or Metadata Repository" on page 2-9

Ultra Search Migration Approaches There are two approaches to migrate user data: the in-place approach and the ETL (extract-transform-load) approach. With the in-place approach, the current `ORACLE_HOME` is used. With the ETL approach, a new `ORACLE_HOME` is created.

Ultra Search In-Place Migration In-place migration upgrades existing configurations and user data to the latest Ultra Search release. Upgraded files are left in place, and the source installation is modified. The benefit to this approach is that it might conserve disk space. With the in-place approach, data migration involves the following six steps:

1. Back up user data
2. Deinstall previous database objects
3. Install new database objects
4. Re-create user instances
5. Restore data
6. Rebuild index

Use the SQL script `wk0upgrade.sql` to run the in-place migration steps one through five, listed in the preceding section. The script is located in the `%ULTRASEARCH_HOME%/admin/` directory. It requires the following input parameters:

- `SYSPW`: password of the user `SYS`
- `WKSYPW`: password of the user `WKSYS`
- `HOST`: database host computer
- `PORT`: database port number
- `ORACLE_SID`: database SID
- `WK_TABLESPACE`: tablespace for Ultra Search
- `WK_TEMPSPACE`: temporary tablespace
- `CONN_STRING`: database connect string
- `ORACLE_HOME`: the path of Oracle home

- `JAVA_EXE_PATH`: Java executable file path
- `PATH_SEPARATOR`: Java classpath separator; use ':' for UNIX or ';' for Windows

The sixth step requires the system administrator to re-activate all crawling schedules through the Ultra Search administration tool.

Ultra Search Extract-Transform-Load Migration Extract-transform-load (ETL) migration extracts the useful subset of configuration data from the source installation, transforms necessary data, and loads or merges this data into a new installation of Ultra Search. This approach might require more disk space, but it offers the following benefits:

- No destabilization of the source installation
- Stability of target installation
- No installer integration requirement

With the ETL approach, data migration involves the following five steps:

1. Install the new system (for example, 9.0.3) in a new `ORACLE_HOME`
2. Re-create user instance schemas and related database objects
3. Re-create user instances
4. Restore data
5. Rebuild index

The first two steps in the ETL approach must be done manually:

- Install Ultra Search 9.0.3 in a separate `ORACLE_HOME`, either on the same computer or on a different computer. If the new 9.0.3 system is installed in the same computer as the old 9.0.1 system, then the database listener port number should be configured to a different number than the old 9.0.1 database. This lets both the old and the new database run at the same time.
- Re-create all Ultra Search 1.0.3 user instance schemas in the new database. Also, for each table data source created in Ultra Search 1.0.3, if the base table is located in the local database, then you must copy the base table to the new database. If the table data source base table is set to a remote database table, then you must re-create the database link from the new database to the remote database.

Use the SQL script `wk0migrate.sql` to run the ETL migration steps three and four. The script is located in the `%ULTRASEARCH_HOME%/admin/` directory. It requires the following input parameters:

- WKSYSWP: password of the user WKSYS
- CONN_STRING: database connect string
- SRC_WKSYSWP: password of the source database (9.0.1 database) user WKSYS
- SRC_CONN_STRING: source database connect string

The fifth step requires the system administrator to re-activate all crawling schedules through the Ultra Search administration tool.

Note: The upgrade script does not roll back the Ultra Search system to the old version if an unexpected error occurs, such as a power failure or system failure.

For in-place migration, back up the database before starting migration. For ETL migration, because all previous data is kept, you can switch back to the previous (for example, 9.0.1) system

Ultra Search Migration Logs The upgrade script provides log files to show which actions the migration has taken. The upgrade script writes the following contents to the log file:

- The current execution step
- Any error message raised from the stored procedures
- Number of data records backup
- Number of data records copied or migrated

For in-place migration, the `wk0upgrade.sql` script writes the execution logs to the file `wk0upgrade.log` in the `%ULTRASEARCH_HOME%/admin/` directory.

For ETL migration, the `wk0migrate.sql` script writes the execution logs to the file `wk0migrate.log` in the `%ULTRASEARCH_HOME%/admin/` directory.

Upgrade from Ultra Search 9.0.2 to 9.0.3

To upgrade Ultra Search 9.0.2 to 9.0.3, perform the following steps:

1. Copy all Ultra Search 9.0.2 files, recursively, under the OracleAS 9.0.2 infrastructure tier `$ORACLE_HOME/ultrasearch/` to a different directory in case if you need to downgrade to 9.0.2 later.
2. Log on to the Ultra Search 9.0.2 administration tool. Stop and disable all crawler synchronization schedules in every Ultra Search instance.

3. Launch Oracle Collaboration Suite release1 installer, and perform the infrastructure install.
4. Specify the directory of the OracleAS 9.0.2 infrastructure as the Oracle Home.
5. The Oracle Universal Installer then detects a previously installed database and automatically upgrades the infrastructure database and the Ultra Search backend (server component).

Upgrade from Ultra Search 9.2 to 9.0.3

Because Ultra Search 9.2 uses the same database schema as Ultra Search 9.0.2, the upgrade procedure is the same.

See Also: "Upgrade from Ultra Search 9.0.2 to 9.0.3" on page 3-14

Configuring the Query Application

The Ultra Search query application is deployed automatically with the Ultra Search installation. However, because Ultra Search allows multiple instances using different schema users, the query application is not configured for how to connect to the database automatically. Database connection is configured by creating a data source in OC4J (not to be confused with an Ultra Search data source). This is done by editing the `data-sources.xml` file.

Step 1: Edit the data-sources.xml File

The `data-sources.xml` file is the OC4J connection management facility. The Ultra Search query application uses OC4J to connect to the database. This is different from the administration tool, because the query user is not a database user; therefore it does not know the database login password.

By editing `data-sources.xml`, the database user and password information is configured with OC4J. The Ultra Search query application finds the data source by using its location, "jdbc/UltraSearchPooledDS".

See Also: "Editing the data-sources.xml File" on page 2-22

Step 2: Deploy Multiple Query Applications Against Multiple Instances

Ultra Search lets multiple instances use different schema users, so multiple query applications can co-exist on the same database.

Each query application requires its database connection information to be defined with `data-sources.xml`. They must be defined to have different location values, such as "jdbc/UltraSearchPooledDS1", "jdbc/UltraSearchPooledDS2", and so on. Correspondingly, the query application must be deployed multiple times in OC4J.

See Also: "Deploying the Sample Query Applications" on page 2-20

Finally, each application deployment must be configured to use the correct entry in `data-sources.xml`. This is done by editing the JSP source for query. For the complete search application, edit `common_customize_instance.jsp` and edit the following line to use the correct location value:

```
String m_datasource_name = "jdbc/UltraSearchPooledDS";
```

Tuning and Performance

This chapter contains the following sections:

- Tuning the Web Crawling Process
- Tuning Query Performance
- Using the Remote Crawler
- Ultra Search on Real Application Clusters
- Table Data Source Synchronization

Tuning the Web Crawling Process

The Ultra Search crawler is a powerful tool for discovering information on Web sites in an organization's intranet. This feature is especially relevant to Web crawling. The other data sources (for example, table or email data sources) are defined such that the crawler does not follow any links to other documents that you might not be aware of.

Web Crawling Strategy

Your Web crawling strategy can be as simple as identifying a few well-known sites that are likely to contain links to most of the other intranet sites in your organization. You could test this by crawling these sites without indexing them. After the initial crawl, you have a good idea of the hosts that exist in your intranet. You could then define separate Web sources to facilitate crawling and indexing on individual sites.

However, in reality, the process of discovering and crawling your organization's intranet is an interactive one characterized by periodic analysis of crawling results and modification to crawling parameters to direct the crawling process somewhat. For example, if you observe that the crawler is spending days crawling one Web host, then you might want to exclude crawling at that host or limit the crawling depth.

Monitoring the Crawling Process

Monitor the crawling process by using a combination of the following methods:

- Monitoring the schedule status with the administration tool
- Monitoring the real time schedule progress with the administration tool
- Monitoring the crawler statistics with the administration tool
- Monitoring the log file for the current schedule

URL Looping

URL looping refers to the scenario where, for some reason, a large number of unique URLs all point to the same document. One particularly difficult situation is where a site contains a large number of pages, and each page contains links to every other page in the site. Ordinarily, this would not be a problem, because the crawler eventually analyzes all documents in the site.

However, some Web servers attach parameters to generated URLs to track information across requests. Such Web servers might generate a large number of unique URLs that all point to the same document.

For example, `http://mycompany.com/somedocument.html?p_origin_page=10` might refer to the same document as `http://mycompany.com/somedocument.html?p_origin_page=13` but the `p_origin_page` parameter is different for each link, because the referring pages are different. If a large number of parameters are specified and if the number of referring links is large, then a single unique document could have thousands or tens of thousands of links referring to it. This is an example of how URL looping can occur.

Monitor the crawler statistics in the Ultra Search administration tool to determine which URLs and Web servers are being crawled the most. If you observe an inordinately large number of URL accesses to a particular site or URL, then you might want to do one of the following:

- **Exclude the Web Server:** This prevents the crawler from crawling any URLs at that host. (You cannot limit the exclusion to a specific port on a host.)
- **Reduce the Crawling Depth:** This limits the number of levels of referred links the crawler will follow. If you are observing URL looping effects on a particular host, then you should take a visual survey of the site to find out an estimate of the depth of the leaf pages at that site. Leaf pages are pages that do not have any links to other pages. As a general guideline, add three to the leaf page depth, and set the crawling depth to this value.

Be sure to restart the crawler after altering any parameters in the Crawler Page. Your changes take effect only after restarting the crawler.

Tuning Query Performance

This section contains suggestions on how to improve the performance of the Ultra Search query. Query performance is generally affected by response time and throughput.

- **Tune the `DB_CACHE_SIZE` initialization parameter.**

The database buffer cache keeps frequently accessed data read from datafiles. Efficient usage of the buffer cache can improve Ultra Search query performance. The cache size is controlled by the `DB_CACHE_SIZE` initialization parameter.

See Also: *Oracle9i Database Performance Tuning Guide and Reference* for information on how to tune this parameter

- Optimize the index.

Optimize the Ultra Search index after the crawler has made substantial updates. To do so, schedule index optimization on a regular basis. Make sure index optimization is scheduled during off-peak hours, because query performance is significantly degraded during index optimization.

See Also: "Index Optimization" on page 7-37

- Optimize the index based on tokens.

Optimize the Ultra Search index by basing it on frequently searched tokens. To log queries, use the administration tool to turn on query statistics collection. The frequently searched tokens then can be passed to `CTX_DDL.OPTIMIZE_INDEX` in token mode. The Ultra Search index name is `WK$DOC_PATH_IDX`.

See Also: *Oracle Text Reference* for more information on `OPTIMIZE_INDEX`

- Simplify query expansion.

The search response time is directly influenced by the Oracle Text query string used. Although Ultra Search provides a default mechanism to expand user input into a Text query, simpler expansions can greatly reduce search time.

See Also:

- "Customizing the Query Syntax Expansion" on page 8-3
- *Oracle Ultra Search API Reference* for the `oracle.ultrasearch.query.Query` interface

- Size the shared pool.

The shared pool stores the library cache and the dictionary cache. The library cache stores recently executed SQL and PL/SQL code. A cache miss on the data dictionary cache or library cache is more expensive than a miss on the buffer cache. For this reason, the shared pool should be sized to ensure that frequently used data is cached. The shared pool size is controlled by the `SHARED_POOL_SIZE` initialization parameter.

See Also: *Oracle9i Database Performance Tuning Guide and Reference* for information on tuning this parameter

- Define JDBC connection pooling.

The Ultra Search middle tier connects to the database through JDBC. Because creation of a connection is an expensive operation in JDBC, a pool of open connections is used to improve the response time of queries. With Oracle Application Server, OC4J can manage the connection pool for the applications.

The minimum size, maximum size, and allocation algorithm of the pool can be specified in the `data-sources.xml` configuration file of OC4J.

The following is an example of a data source definition, with minimum 2 and maximum 30 open-connections. Each connection closes after 30 seconds of inactivity, and new connections are created dynamically according to load. The other caching schemes are `FIXED_WAIT_SCHEME` and `FIXED_RETURN_NULL_SCHEME`.

Note: `DYNAMIC_SCHEME = 1`, `FIXED_WAIT_SCHEME = 2`, and `FIXED_RETURN_NULL_SCHEME = 3`

```
<data-source
  class="oracle.jdbc.pool.OracleConnectionCacheImpl"
  name="UltraSearchDS"
  location="jdbc/UltraSearchPooledDS"
  username="user"
  password="pass"
  url="jdbc:oracle:thin:@hostname:1521:oracle_sid"
  min-connections="2"
  max-connections="30"
  inactivity-timeout="30" >
  <property name="cacheScheme" value="1" />
</data-source>
```

- Pin the query package in memory.

Pin frequently used packages in the shared memory pool. When a package is pinned, it remains in memory, no matter how full the pool gets or how frequently you access the package. You can pin packages using the supplied package `DBMS_SHARED_POOL`.

The PL/SQL package used for Ultra Search query is `WKSYS.WK_QRY`.

See Also: *Oracle9i Supplied PL/SQL Packages and Types Reference*

Using the Remote Crawler

Without the Ultra Search remote crawler, you must run the Ultra Search crawler on the same host as the Oracle Database. For large data sets, you can improve performance by running the Ultra Search crawler on one or more separate hosts from the Oracle Database. Because the Ultra Search crawler is a pure Java application, it communicates with the Oracle Database through JDBC.

Ultra Search remote crawler instances are always launched by the Oracle Database. By launching remote crawlers from the Oracle Database, you can leverage the high-availability of the Oracle Database. The Ultra Search scheduling mechanism runs within the Oracle Database and therefore automatically uses the database's high availability features.

The Oracle Database uses Java remote method invocation (RMI) to communicate with the remote crawler hosts. Therefore, each remote host must have an RMI registry and an RMI daemon running.

1. When a crawling schedule is activated, the Ultra Search scheduler launches a Java program as a separate process on the database host. This Java program is known as the `ActivationClient`.
2. This program attempts to connect to the remote crawler host through the standard RMI registry and RMI daemon on ports 1098 and 1099. If successful, then the `ActivationClient` receives a remote reference to a Java object running on the remote host. This remote Java object is known as the `ActivatableCrawlerLauncher`.
3. The `ActivationClient` then instructs the `ActivatableCrawlerLauncher` to launch the Ultra Search crawler on the remote host. The `ActivatableCrawlerLauncher` launches the Ultra Search crawler as a separate Java process on the remote host.

Note: By default, RMI sends data over the network unencrypted. Using the remote crawler to perform crawling introduces a potential security risk. A malicious entity within the enterprise could steal the Ultra Search instance schema and password by listening to packets going across the network. Refrain from using the remote crawler feature if this security risk is unacceptable.

Scalability and Load Balancing

Each Ultra Search schedule can be associated with exactly one crawler. The crawler can run locally on the Oracle database host or on a remote host. There is no limit to the number of schedules that can be run. Similarly, there is no limit to the number of remote crawler hosts that can be run. However, each remote crawler host requires that the Ultra Search middle tier be installed on its host.

By using several remote crawler hosts and carefully allocating schedules to specific hosts, you can achieve scalability and load balancing of the entire crawling process.

Installation and Configuration Sequence

1. Make sure that you have installed the Ultra Search backend (server component) on the Oracle Database, the Ultra Search middle tier on one or more Web server hosts, and the Ultra Search middle tier on all remote crawler hosts.

See Also: Chapter 2, "Installing and Configuring Ultra Search"

2. Export the following common resources on the database host:
 - The temporary directory
 - The log directory
 - The mail archive directory (if you are using the Ultra Search mailing list feature)

These resources are merely directories that must be accessible by all remote crawler instances over the network. Use whatever mechanism you want to share these resources with a remote crawler host.

The remote crawler code is pure Java. Therefore, it is platform-independent. For example, your Ultra Search installation might consist of four hosts: one database server (host X) running Solaris on which the Ultra Search backend (server component) is installed; one remote crawler host (host Y1) running on Windows; one remote crawler host (host Y2) running on Solaris; and one remote crawler host (host Y3) running on Linux.

In this scenario, export the shared directories on host X using the UNIX "export" command. Then use the UNIX "mount" command on hosts Y2 and Y3 to mount the exported directories. For host Y1, you must purchase a third-party NFS client for Windows and use that to mount the shared directories. If host X is a

Linux server, then you can create Samba shares and thereby mount those shares on Windows without needing any third party software.

3. Configure the remote crawler with the administration tool.

Edit that profile by manually entering all mount points for the shared crawler resources that you defined. To edit the remote crawler profile, navigate to the **Crawler: Remote Crawler Profiles** page and click **Edit** for the remote crawler profile you want to edit. Specify values for the following parameters:

- Mount point for temporary directory path as seen by the remote crawler
- Mount point for log directory path as seen by the remote crawler
- Mount point for mail archive path as seen by the remote crawler (if you are using the Ultra Search mailing list feature)

Additionally, you must specify the following crawler parameters before you can begin crawling:

- Number of crawler threads that the remote crawler uses for gathering documents
- Number of processors on the remote crawler host

4. Complete the crawler configuration with the administration tool.

The minimum set of parameters that likely need to be configured are the following:

- Seed URLs
- Web proxy
- A schedule

Each schedule must be assigned to a remote crawler or the local crawler. (The local crawler is the crawler that runs on the local Oracle database host itself). To assign the a schedule to a remote crawler host or the local database host, click the host name of a schedule in the **Schedules** page.

You can also turn off the remote crawler feature for each schedule, thereby forcing the schedule to launch a crawler on the local database host, instead of the specified remote crawler host. To turn off the remote crawler feature, click the host name of a schedule in the **Synchronization Schedules** page. If a remote crawler host is selected, then you can enable or disable the remote crawler.

See Also: Chapter 7, "Understanding the Ultra Search Administration Tool"

5. Start up the RMI registry and RMI daemon on each remote crawler host.

Use the helper scripts in `$ORACLE_HOME/tools/remotecrawler/scripts/operating_system` to do this.

- If the remote crawler is running on a UNIX platform, then source the `$ORACLE_HOME/tools/remotecrawler/scripts/unix/runall.sh` Bourne shell script.
- If the remote crawler is running on a Windows host, then run the `%ORACLE_HOME%\tools\remotecrawler\scripts\winnt\runall.bat` file.

The `runall.sh` and `runall.bat` scripts perform the following tasks in sequence:

- `define_env` is invoked to define necessary environment variables
- `runregistry` is invoked to start up the RMI registry
- `runrmid` is invoked to start up the RMI daemon
- `register_stub` is invoked to register the necessary Java classes with the RMI subsystem

You can invoke `runregistry`, `runrmid`, and `register_stub` individually. However, you must first invoke `define_env` to define the necessary environment variables.

6. Launch the remote crawler from the administration tool, and verify that it is running.

The state of the schedule is listed in the **Schedules** page. The remote crawler launching process takes up to 90 seconds to change state from `LAUNCHING` to `FAILED`, if failure occurs.

To view the schedule status, click the crawler status in the schedules list. To view more details, especially in the event of failure, click the schedule status itself. This brings up a detailed schedule status.

The remote crawler fails to launch if any one of the following requirements are not met:

- The RMI registry is not running and listening on port 1099 of each remote host.
- The RMI daemon is not running and listening on port 1098 of each remote host.

- The necessary Java objects have not been successfully registered with each RMI registry.

After a remote crawler is launched, verify that it is running by one or more of the following methods:

- Check for active Java processes on the remote crawler host.

A simple way to confirm that remote crawler is running on the remote crawler host is to use an operating system command, such as `ps` on UNIX systems. You should look for active Java processes.

- Monitor the contents of the schedule log file.

If the remote crawler is running successfully, you should see the contents of the schedule log file changing periodically. The schedule log file is located in the shared log directory.

Ultra Search on Real Application Clusters

Ultra Search can crawl on one fixed node or on any node, depending on the storage access configuration of the Real Application Clusters system. PL/SQL APIs are provided to specify which node should run the crawler, if needed. For Ultra Search administration and the Ultra Search query application, you can configure the connection string to connect to any node of Real Application Clusters.

See Also: The documentation for Oracle Database Real Application Clusters

Configuring Storage Access

The disk of any node in a Real Application Clusters system can be shared (cluster file system) or not shared (raw disk). For Real Application Clusters on a cluster file system (CFS), the cache files generated by the crawler on any node are visible to any Oracle instance and can be indexed by any Oracle instance that performs index synchronization. If the disk is not shared, then the crawler must run on one particular Oracle instance to ensure that all cache files can be indexed.

This is due to the nature of Oracle Text indexing, where rows inserted into one table by different sessions go to the same pending queue, and whoever initiates index synchronization attempts to index all of the inserted rows. Because of this limitation, on a CFS, Ultra Search is configured to launch the crawler on any database instance. If it is not on a CFS, then Ultra Search launches the crawler on the database instance where `INSTANCE_NUMBER = 1`.

The Ultra Search administrator can configure which instance runs the crawler with the following PL/SQL API:

```
WK_ADM.SET_LAUNCH_INSTANCE(instance_name, connect_url);
```

where `instance_name` is the name of the launching instance (or the database name if it is to be launched on any node) and `connect_url` is the connect descriptor.

For connection to a single database instance, the descriptor can be in the short form "*host:port:SID*" or the connect descriptor (Oracle Net keyword-value pair). For example:

```
(DESCRIPTION= (ADDRESS_
LIST= (ADDRESS= (PROTOCOL=TCP) (HOST=c1s02a) (PORT=3999) ) ) (CONNECT_DATA= (
SERVICE_NAME=acme.us.com) ) )
```

To connect to any database instance, the full database connect descriptor must be used. For example:

```
(DESCRIPTION= (LOAD_BALANCE=yes) (ADDRESS_
LIST= (ADDRESS= (PROTOCOL=TCP) (HOST=c1s02a) (PORT=3999
) ) (ADDRESS= (PROTOCOL=TCP) (HOST=c1s02b) (PORT=3999) ) ) (CONNECT_DATA= (SERVICE_
NAME=acme.us.com) ) )
```

See Also: *Oracle9i JDBC Developer's Guide and Reference* for configuration details.

You cannot configure Ultra Search to launch the crawler on any node on a non-cluster file system.

To query on the existing launching instance configuration, use the following PL/SQL API:

```
WK_ADM.GET_LAUNCH_INSTANCE RETURN VARCHAR2;
```

This returns the name of the launching instance or the database name if any node can launch the crawler.

Remote Crawler File Cache

The Ultra Search remote crawler requires that the remote file system be mounted on the Oracle instance for indexing.

For cluster file system Real Application Clusters, the file system of the remote computer should be NFS mounted to all nodes of the system.

For non-cluster file system Real Application Clusters, the NFS mount can be limited to the specific node where the Oracle instance is serving the remote crawler. There is no advantage to mounting the remote file system to all nodes--it could lead to stale NFS handles when nodes go down. When there is a configuration change to move to a different Oracle instance, the remote file system should be NFS mounted to the new node accordingly.

Logging on to the Oracle Instance

All components of Ultra Search use the JDBC Thin Driver with the connect string consisting of "*hostname:port:SID*" or the full connect descriptor as seen in `tnsnames.ora`.

The administration middle tier connects to the Oracle database with a JDBC connection specified in the `ultrasearch.properties` file. If the client serving node is down, then you must manually edit the `ultrasearch.properties` file to connect to a different Oracle instance.

Query Search Application for Real Application Clusters

Query components should fully utilize Real Application Clusters. You can specify the JDBC connection string as a database connect descriptor so that it can connect to any Oracle instance in Real Application Clusters. For example:

```
"jdbc:oracle:thin:@(DESCRIPTION=(LOAD_BALANCE=yes) (ADDRESS_
LIST=(ADDRESS=(PROTOCOL=TCP) (HOST=c1s02a) (PORT=3999
)) (ADDRESS=(PROTOCOL=TCP) (HOST=c1s02b) (PORT=3999))) (CONNECT_DATA=(SERVICE_
NAME=acme.us.com))) "
```

See Also: *Oracle9i JDBC Developer's Guide and Reference*

Java Crawler

The connect string used by Ultra Search crawler is initialized during installation and can be changed with the `WK_ADM.SET_LAUNCH_INSTANCE` API. When there is a system configuration change, such as adding or dropping a node, the connect string is changed automatically.

Choosing a JDBC Driver

The Ultra Search administrator optionally can configure the local crawler to use the JDBC OCI driver to log on to the database. This is done with the following PL/SQL API:

```
WK_ADM.SET_JDBC_DRIVER(driver_type)
```

Where

- Thin driver (default) `driver_type = 0`
- OCI driver `driver_type = 1`

This API requires super-user privileges. The change affects all Ultra Search instances.

Note: The OCI driver requires that environment variables, such as `LD_LIBRARY_PATH` and `NLS_LANG`, be set properly on the launching database instance. The crawler inherits the environment setting from the Oracle process. Therefore, you must configure them appropriately before starting Oracle.

See Also: *Oracle9i JDBC Developer's Guide and Reference* for configuration details on using the OCI driver.

The following PL/SQL API determines which kind of JDBC drivers are used currently:

```
WK_ADM.GET_JDBC_DRIVER RETURN NUMBER;
```

Table Data Source Synchronization

Ultra Search crawls database tables in the local Oracle Database instance where Ultra Search is installed. Additionally, it can crawl remote databases if they have been linked to the main Oracle Database. Remote databases are linked to the main Oracle instance with database links.

See Also: *Oracle9i Database Administrator's Guide* for instructions on how to create database links

Ultra Search provides a logging mechanism to optimize crawling of table sources. Using this logging mechanism, only newly updated documents are revisited during the crawling process. If the source database is not an Oracle database, then you must perform a sequence of steps to use this feature.

Synchronizing Crawling of Oracle Databases

Before creating log tables and log triggers, make sure that the Ultra Search instance schema has the `CREATE ANY TABLE` and `CREATE ANY TRIGGER` system privileges. For tables in Oracle databases, data definition language (DDL) statements are provided to create the following:

Create Log Table

The log table stores changes that have occurred in the base table. The Ultra Search crawler uses the change information to figure out which rows need to be recrawled. For example, a log table generated by Ultra Search could be named `WK$LOG`.

The structure of the log table conforms to the following rules:

- For every primary key column of the base table, a column must be created in the log table.
- There can be up to only eight primary key columns in the base table.
- Each column in the log table that corresponds to a primary key column must be named `Kx`, where `x` is a number from one to eight.
- Each column in the log table that corresponds to a primary key column must be of type `VARCHAR2(1000)`.
- There must be exactly one column named `mark` that has type `CHAR(1)`.
- The column named `mark` must have a default value `F`.

For example, the base table `employees` has the following structure:

Column Name	Column Type
ID	NUMBER
NAME	VARCHAR2(200)
ADDRESS	VARCHAR2(400)
TELEPHONE	VARCHAR2(10)
USERNAME	VARCHAR2(24)

If the primary key of the `employees` table comprises of the `ID` and `NAME` columns, then a log table `WK$LOG` (whose name is generated automatically) is created with the following structure:

Column Name	Column Type
K1	NUMBER
K2	VARCHAR2 (200)

The SQL statement for creating the log table is as follows:

```
CREATE TABLE WK$LOG (
K1 VARCHAR2(1000),
K2 VARCHAR2(1000),
MARK CHAR(1) default 'F')
```

Create Log Triggers

An INSERT trigger, UPDATE trigger, and DELETE trigger are created. The Oracle trigger definitions are as follows:

INSERT Trigger Statement Every time a row is inserted into the `employees` base table, the INSERT trigger inserts a row into the log table. The row in the log table records the new values of the `id` and the `name` into the `k1` and `k2` columns. An F is inserted into the `mark` column to signal the crawler that work needs to be done for this row.

For example:

```
CREATE OR REPLACE TRIGGER wk$ins
AFTER INSERT ON employees
FOR EACH ROW;

BEGIN
  INSERT INTO WK$LOG(k1,k2,mark)
    VALUES (:new.id, :new.name, 'F');
END;
```

UPDATE Trigger Statement Every time a row is updated in the `employees` base table, the UPDATE trigger inserts two rows into the log table. The first row in the log table records the old values of the `id` and the `name` into the `k1` and `k2` columns. An F is inserted into the `mark` column to signal the crawler that work needs to be done for this row. The second row in the log table records the new values of the `id` and the `name` into the `k1` and `k2` columns.

For example:

```
CREATE OR REPLACE TRIGGER wk$upd
AFTER UPDATE ON employees
FOR EACH ROW;

BEGIN
    INSERT INTO WK$LOG(k1,k2,mark)
        VALUES (:old.id, :old.name, 'F');
    INSERT INTO WK$LOG(k1,mark)
        VALUES (:new.id, :new.name, 'F');
END;
```

DELETE Trigger Every time a row is deleted from the `employees` base table, the `DELETE` trigger inserts a row into the log table. The row in the log table records the old values of the `id` and the `name` into the `k1` and `k2` columns. An `F` is inserted into the `mark` column to signal the crawler that work needs to be done for this row.

For example:

```
CREATE OR REPLACE TRIGGER wk$del
AFTER DELETE ON employees
FOR EACH ROW;

BEGIN
    INSERT INTO WK$LOG(k1,k2,mark)
        VALUES (:old.id, :old.name, 'F');
END;
```

Synchronizing Crawling of Non-Oracle Databases

For tables in non-Oracle remote databases, you must perform the following steps:

1. Manually create the log table. The log table must conform to the rules for log tables described earlier. Also, it must reside in the same schema and database instance as the base table.
2. Create three triggers that record inserts, updates, and deletes on the base table. These triggers must exhibit the same behavior as the triggers described earlier for Oracle tables.
3. Associate the log table. When you have completed these tasks, choose the "Enable logging mechanism (non-Oracle tables)" option during the creation of an Ultra Search table data source. By choosing that option, the Ultra Search administration tool prompts you for the name of the log table in the remote database. Ultra Search associates this log table with the base table. Ultra Search assumes that you have correctly performed steps 1 and 2.

Security in Ultra Search

The ability to control user access to Web content is critical. This chapter describes the architecture and configuration of security for Ultra Search.

This chapter contains the following sections:

- About Ultra Search Security
- Configuring a Security Framework for Ultra Search
- Configuring Ultra Search Security

See Also:

- *Oracle Application Server 10g Security Guide* for an overview of Oracle Application Server security and its core functionality
- *Oracle Identity Management Concepts and Deployment Planning Guide* for guidance on the Oracle security infrastructure

About Ultra Search Security

This section describes the Ultra Search security model. It contains the following sections:

- Ultra Search Security Model
- Classes of Users and Their Privileges
- Ultra Search Admin Privilege Model in the Hosted Environment
- Resources Protected by Ultra Search
- Authorization and Access Enforcement
- How Ultra Search Leverages Security Services
- How Ultra Search Leverages the Identity Management Infrastructure
- Ultra Search Extensibility and Security

Ultra Search Security Model

Security problems, such as unauthorized access to information, can lead to loss of productivity. Search engines like Oracle Ultra Search provide access to a vast variety of content repositories in a single gateway. Each one of these repositories has its own security model that determines whether a particular end user can access a particular document. Because Ultra Search provides access to data from multiple repositories, existing security information in each repository must be carefully supported to avoid unauthorized access.

This section describes the security architecture of Ultra Search. Security is implemented at the following levels:

- User authentication
This is the identification of a user, through LDAP and Oracle Internet Directory, at Ultra Search front-end interfaces.
- User entitlement
This determines whether a user can access information about a particular item in the results list. It is implemented by access control lists (ACLs). Ultra Search provides mapped-security to third-party repositories by retrieving the access control list for each document at the time of indexing and storing them in Ultra Search. Ultra Search does not need any connection with the repository itself to validate access privileges.

- **Secure communications**
All content crawling, indexing, and querying is encrypted using secure socket layer (SSL), a worldwide standard for encryption over the HTTP protocol (HTTPS).
- **Security of Ultra Search**
Actual Ultra Search security is handled by the dictionary data in the Ultra Search database, the administrative user, and password data.

Classes of Users and Their Privileges

To grant an Ultra Search user administration privileges, you must assign the user to an administration group. Each user can belong to one or more groups. The following groups are created for each Ultra Search instance:

- **Instance administrators:** Users in this group can only manage instances for which they have privileges.
- **Super-users:** Users in this group can manage all instances, including creating instances, dropping instances, and granting privileges.

Ultra Search also has two classes of users:

- **Single Sign-on (SSO) users:** These users are managed by the Oracle Internet Directory (OID) and are authenticated by the SSO server. The Ultra Search administration tool identifies all Ultra Search instances to which the SSO user has access. This is available only if you have the Oracle Identity Management infrastructure installed.
- **Database users (non-SSO):** These users exist in the database on which Ultra Search runs.

Ultra Search Default Users

New Ultra Search instances contain the following users:

- **WK_TEST:** This is the instance administrator user that hosts the default instance, called WK_INST. In other words, WK_TEST is the instance administrator for WK_INST. For security purposes, WK_TEST is locked after the installation. The administrator should login to the database as DBA role, unlock the WK_TEST user account, and set the password to be WK_TEST. (The password expires after the installation.) If the password is changed to anything other than WK_TEST, then you must also update the cached schema password using the

administration tool **Edit Instance** page after you change the password in the database.

- **WKSYS:** This is a database super-user. WKSYS can grant super-user privileges to other users, such as WK_TEST. All Ultra Search database objects are installed in the WKSYS schema.

Note: The WKUSER role is required to host instances.

Ultra Search Admin Privilege Model in the Hosted Environment

In a hosted environment, one enterprise (for example, an application service provider) makes Ultra Search available to other enterprises and stores information for them. The enterprise performing the hosting is called the **default subscriber**, and the enterprises that are hosted are called **subscribers**.

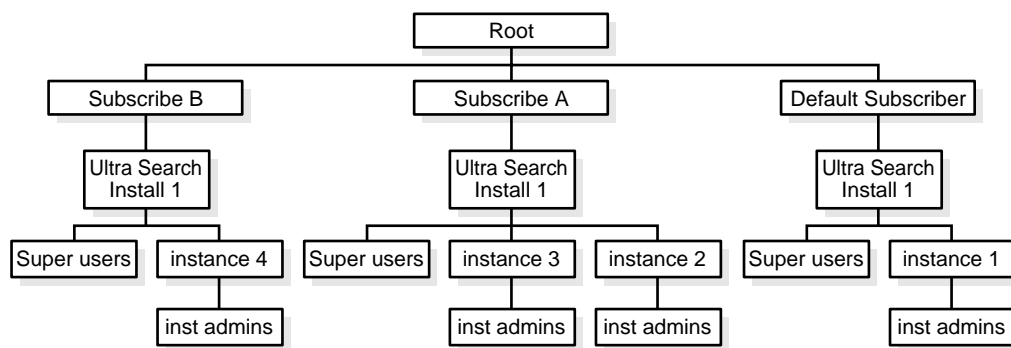
Note: This is available with the Oracle Application Server release and the Oracle Collaboration Suite release. This is not available with the Oracle Database release.

The default subscriber and its search base are specified in the following attributes of the OID entry "cn=Common,cn=Products,cn=OracleContext":

- orclDefaultSubscriber
- orclSubscriberSearchBase

In a non-hosted environment, in which there are no subscribers, the enterprise installing Ultra Search is the default subscriber. All Ultra Search administration groups (super-user and instance administrator groups) are created under Default Subscriber Oracle Context (for example, cn=OracleContext,dc=us, dc=oracle, dc=com) in the OID Directory Information Tree (DIT).

Figure 5-1 shows an example of the OID topology of a hosted environment. There are two subscribers (A and B) and the default subscriber. Each subscriber has its own super-user privilege group associated with it. There are four Ultra Search instances created in the Ultra Search back-end 'install 1'. 'Instance 1' is associated with the default subscriber. 'Instance 2' and 'Instance 3' are associated with 'Subscriber A'. 'Instance 4' is associated with 'Subscriber B'. Each Ultra Search instance has its instance administration group associated with it.

Figure 5–1 *OID Topology of a Hosted Environment*

Admin Privilege Model

This section describes the privilege model of Ultra Search administration tool in the hosted environment. The model applies to both the SSO login mode and the non-SSO login mode.

Privileges to Administer an Ultra Search Instance In non-SSO mode, only database users can login to the admin tool.

- If the login database user has the super-user privilege, then the user can administer all Ultra Search instances across the default subscriber and any other subscribers.
- If the login database user only has the admin privilege on a particular Ultra Search instance, then the user can administer the instance regardless of whether the instance is associated with the default subscriber or any other subscribers.

In SSO mode, only SSO users can login to the admin tool

- If the SSO user belongs to the default subscriber, then the following is true:
 - If the SSO user has the super-user privilege, then the user can administer all Ultra Search instances across the default subscriber and any other subscribers (for example, instances 1,2,3,4).
 - If the SSO user has the admin privilege on a particular Ultra Search instance (for example, instance1) within the default subscriber, then the user can administer the instance (instance 1) that is associated with the default subscriber.
- If the SSO user belongs to a subscriber, then the following is true:

- If the SSO user has the super-user privilege, then the user can administer only Ultra Search instances within the subscriber to which he belongs. (For example, if the user from subscriber A has the super-user privilege, then the user can only administer instance 2 and 3.)
- If the SSO user has the admin privilege on a particular Ultra Search instance (for example, instance 2), then the user can administer the instance (instance 2) that is associated with the subscriber (subscriber A).

Privileges to Create and Drop an Ultra Search Instance To create or drop an Ultra Search instance, the user (either the database or the SSO user) must have the super-user privilege.

In non-SSO mode, the database user can create or drop an instance and associate the instance with any subscriber, including the default subscriber.

In SSO mode:

- If the login SSO user belongs to the default subscriber, then the user can create or drop an instance and associate the instance with any subscriber, including the default subscriber.
- If the login SSO user belongs to a particular subscriber, then when the user creates an instance, the instance is created and associated with the subscriber to which the login user belongs. Because the user might not have access to create the instance schema in the database, the user must inform the hosting company (default subscriber) to create the database schema for hosting the instance.

Privileges to Grant or Revoke a Super-User To grant or revoke a super-user, login to the administration tool as a super-user.

In non-SSO mode (database user login), only super-users can grant or revoke the super-user privilege to and from other database users.

In SSO mode:

- If the login SSO user belongs to the default subscriber, then the user can do the following:
 - Grant or revoke the super-user privilege to SSO users in the default subscriber.
 - Grant or revoke the super-user privilege to SSO users in a particular subscriber.

- If the login SSO user belongs to a particular subscriber, then the user can grant or revoke the super-user privilege to users within the same subscriber to which the login user belongs.

Privileges to Grant or Revoke an Instance Administrator To grant or revoke an instance administrator, login to the admin tool as a super-user or an instance administrator.

In non-SSO mode (database user login), only super-users or instance administrators can grant or revoke the instance admin privilege to and from other database users.

In SSO mode:

- The login SSO user can grant or revoke only the instance admin privilege to SSO users within the subscriber the instance as associated with. For example, the user can grant the admin privilege on 'Instance 2' or 'Instance 3' to an SSO user in subscriber A.
- The login SSO user cannot grant or revoke the instance admin privilege to SSO users within a different subscriber. For example, the user cannot grant the admin privilege on 'Instance 2' or 'Instance3' to an SSO user in subscriber B.

Resources Protected by Ultra Search

All publicly crawled data is publicly accessible.

The following resources are protected by Ultra Search:

- Access control list (ACL)-aware crawled data is protected; in other words, it is private to users named by the ACL.
- All passwords are protected.
- User-defined data source parameters are protected.

Authorization and Access Enforcement

There are three possible entry points to Ultra Search:

- The database: This contains all data. All data and metadata is protected with row level security. All passwords are encrypted.
- The Ultra Search administration tool: This does not contain crawled data. You must authenticate with SSO or database authentication.
- The Ultra Search query tool: This contains crawled data. Unauthenticated users can see only public data. Authenticated users can see public data and

ACL-protected information. Users must authenticate themselves to see private information.

How Ultra Search Leverages Security Services

Ultra Search uses the following to leverage security services:

- Ultra Search uses secure socket layers (SSL), the industry standard protocol for managing the security of message transmission on the Internet. This is used for securing RMI connections, HTTPS crawling, and secure JDBC.
- JAZN: OracleAS Containers for J2EE (OC4J) implements a Java authentication and authorization service (JAAS) provider called JAZN. This provides application developers with user authentication, authorization, and delegation services to integrate into their application environments.

See Also: "Configure a Secure Ultra Search Installation" on page 2-6

How Ultra Search Leverages the Identity Management Infrastructure

Ultra Search uses the SSO server and OID to leverage the Oracle Identity Management infrastructure.

With the SSO server, you can log on once for all components, and the Ultra Search administrative interface allows user management operations on either database users or SSO users. Authenticated SSO users never see the Ultra Search logon screen. Instead, they can immediately choose an instance. The Ultra Search administration tool and the query tool use SSO.

Oracle Internet Directory (OID) is Oracle's native LDAP v3-compliant directory service, built as an application on top of the Oracle database. OID hosts the Oracle common identity. All Ultra Search instances are registered with OID.

See Also: "Integration with Oracle Internet Directory" on page 1-11

Ultra Search has native identity management; therefore, in the absence of the identity management infrastructure, Ultra Search uses native user management available with the Oracle database.

Ultra Search Extensibility and Security

Ultra Search is extensible (for example, the crawler agent), but this poses no extra security considerations.

Configuring a Security Framework for Ultra Search

This section describes special security configuration steps within Ultra Search.

Configuring Security Framework Options for Ultra Search

Storing clear text passwords in `data-sources.xml` poses a security risk. Avoid this by using password indirection to specify the password. This lets you enter the password in `jazn-data.xml`, which is automatically encrypted, and point to it from `data-sources.xml`.

See Also:

- "Editing the `data-sources.xml` File" on page 2-22
- *Oracle Application Server Containers for J2EE Services Guide*

Configuring Oracle Identity Management Options for Ultra Search

To configure the Ultra Search administration tool with the SSO server, you must follow certain steps.

See Also: "Configuring the Administration Tool with Single Sign-On Server" on page 2-18

Configuring Ultra Search Security

Ultra Search has no specific security passwords.

See Also: "Configuring Security Framework Options for Ultra Search" on page 5-9 for more information on Ultra Search configuration issues to leverage security

Understanding the Ultra Search Crawler and Data Sources

This chapter contains the following topics:

- Overview of the Ultra Search Crawler
- Crawler Settings
- Crawler Data Sources
- Document Attributes
- Crawling Process for the Schedule
- Data Synchronization
- Ultra Search Remote Crawler

See Also: "Tuning Query Performance" on page 4-3

Overview of the Ultra Search Crawler

The Ultra Search crawler is a Java process activated by your Oracle server according to a set schedule. When activated, the crawler spawns processor threads that fetch documents from various data sources. These documents are cached in the local file system. When the cache is full, the crawler indexes the cached files using Oracle Text. This index is used for querying.

Note: An empty index is created when an Ultra Search instance is created. You can alter the index using SQL. The existing preferences, such as language-specific parameters, are defined in the `$ORACLE_HOME/ultrasearch/admin/wk0pref.sql` file.

Crawler Settings

Before you can use the crawler, you must set its operating parameters, such as the number of crawler threads, the crawler timeout threshold, the database connect string, and the default character set. To do so, use the **Crawler Settings Page** in the administration tool.

See Also: "Crawler Page" on page 7-12

Crawler Data Sources

In addition to the Web access parameters, you can define specific data sources on the **Sources** page in the administration tool. You can define one or more of the following data sources:

- Web sites
- Database tables
- Files
- Mailing lists
- OracleAS Portal page groups
- User-defined data sources (requires crawler agent)

Using Crawler Agents

If you are defining a user-defined data source to crawl and index a proprietary document repository or management system, such as Lotus Notes or Documentum, then you must implement a crawler agent as a Java class. The agent collects document URLs and associated metadata from the proprietary document source and returns the information to the Ultra Search crawler, which enqueues it for later crawling. For more information on defining a new data source type, see the User-Defined sub-tab in **Sources** page in the administration tool.

Synchronizing Data Sources

You can create synchronization schedules with one or more data sources attached to it. Synchronization schedules define the frequency at which the Ultra Search index is kept up to date with existing information in the associated data sources. To define a synchronization schedule, use the **Sources** page in the administration tool.

Display URL and Access URL

For some applications, for security reasons, the URL crawled is different from the one seen by the end user. For example, crawling on an internal Web site inside a firewall might be done without security checking, but when queried by the end user, a corresponding mirror URL outside the firewall must be used. This mirror URL is called the display URL.

By default, the display URL is treated as the access URL unless a separate access URL is provided. The display URL must be unique in a data source; so two different access URLs cannot have the same display URL.

See Also: "Sources Page" on page 7-20

Document Attributes

Document attributes, or metadata, describe the properties of a document. Each data source has its own set of document attributes. The value is retrieved during the crawling process and then mapped to one of the search attributes and stored and indexed in the database. This lets you query documents based on their attributes. Document attributes in different data sources can be mapped to the same search attribute. Therefore, you can query documents from multiple data sources based on the same search attribute.

If the document is a Web page, the attribute can come from the HTTP header or it can be embedded inside the HTML in metatags. Document attributes can be used

for many things, including document management, access control, or version control. Different data sources can have attributes of different names which are used for the same purpose; for example, "version" and "revision". It can also have the same attribute name for different purposes; for example, "language" as in natural language in one data source but as programming language in another.

Search attributes are created in three ways:

- System-defined search attributes, such as title, author, description, subject, and mimetype
- Search attributes created by the system administrator
- Search attributes created by the crawler. (During crawling, the crawler agent maps the document attribute to a search attribute with the same name and data type. If not found, then the crawler creates a new search attribute with the same name and type as the document attribute defined in the crawler agent.)

The list of values (LOV) for a search attribute can help you specify a search query. If attribute LOV is available, then the crawler registers the LOV definition, which includes attribute value, attribute value display name, and its translation.

Crawling Process for the Schedule

The first time the crawler runs, it must fetch Web pages, table rows, files, and so on based on the data source. It then adds the document to the Ultra Search index. The crawling process for the schedule is broken into two phases:

1. Queuing and Caching Documents
2. Indexing Documents

Queuing and Caching Documents

Figure 6-1 on page 6-6 and Figure 6-2 on page 6-7 illustrate an instance of the crawling cycle in a sequence of nine steps. The example uses a Web data source, although the crawler can also crawl other data source types.

Figure 6-1 illustrates how the crawler and its crawling threads are activated. It also shows how the crawler queues hypertext links to control its navigation. This figure corresponds to Steps 1 to 5.

Figure 6-2 illustrates how the crawler caches Web pages. This figure correspond to Steps 6 to 8.

The steps are the following:

1. Oracle spawns the crawler according to the schedule you specify with the administration tool. When crawling is initiated for the first time, the URL queue is populated with the seed URLs. Figure 6-1.
2. Crawler initiates multiple crawling threads.
3. Crawler thread removes the next URL in the queue.
4. Crawler thread fetches the document from the Web. The document is usually an HTML file containing text and hypertext links.
5. Crawler thread scans the HTML file for hypertext links and inserts new links into the URL queue. Duplicate links already in the document table are discarded.
6. Crawler caches the HTML file in the local file system. Figure 6-2 on page 6-7.
7. Crawler registers URL in the document table.
8. Crawler thread starts over by repeating Step 3.

Fetching a document, as shown in Step 4, can be time-consuming because of network traffic or slow Web sites. For maximum throughput, multiple threads fetch pages at any given time.

Note: URLs remain visible until the next crawling run. When the crawler detects that the URL is no longer there, it is removed from the `wk$doc` table where Oracle Text automatically marks this document as deleted, even though the index data still exists. Cleanup is done through index optimization, which can be scheduled separately.

Figure 6–1 Queuing URLs

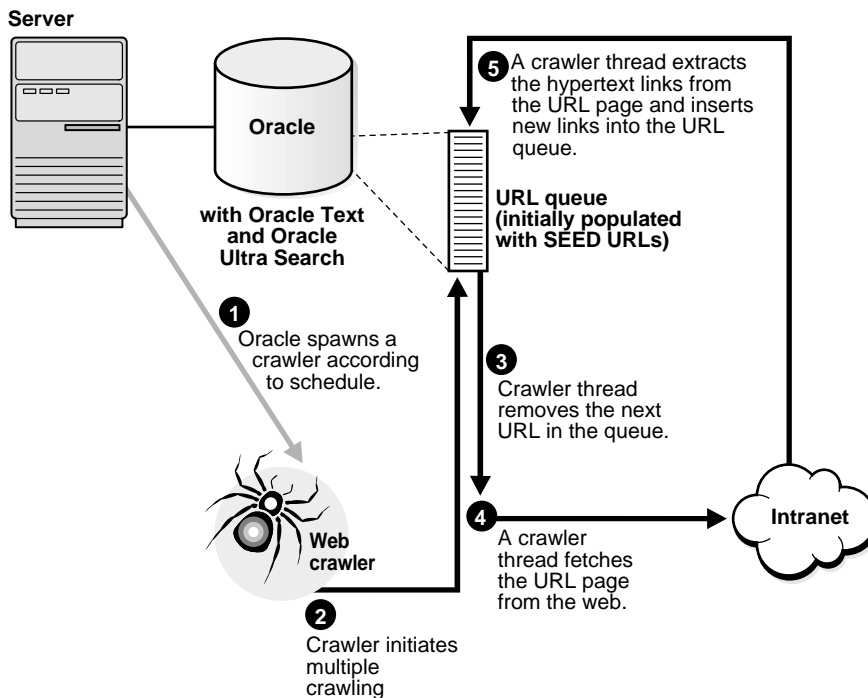
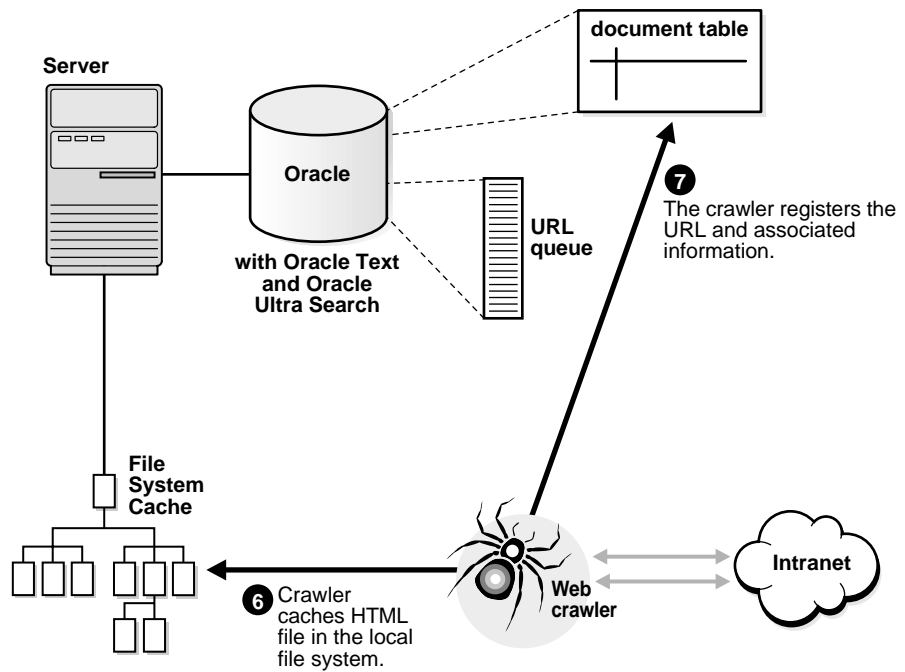
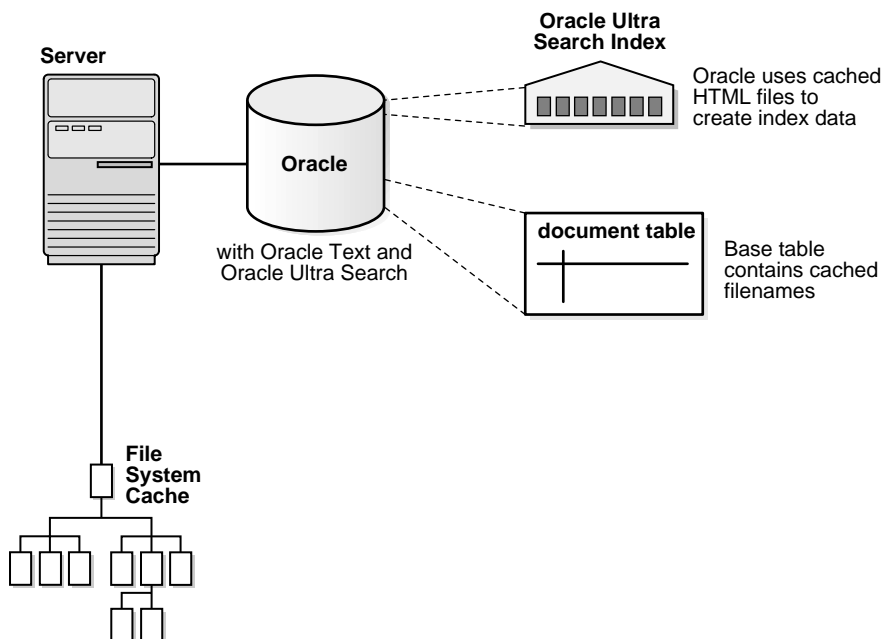


Figure 6-2 Caching URLs



Indexing Documents

When the file system cache is full (default maximum size is 20 megabytes), document caching stops and indexing begins. In this phase, Ultra Search augments the Oracle Text index using the cached files referred to by the document table. See Figure 6-3.

Figure 6–3 Indexing Documents

Data Synchronization

After the initial crawl, a URL page is only crawled and indexed if it has changed since the last crawl. The crawler determines if it has changed with the HTTP `If-Modified-Since` header field or with the checksum of the page. URLs that no longer exist are marked and removed from the index.

To update changed documents, the crawler uses an internal checksum to compare new Web pages with cached Web pages. Changed Web pages are cached and marked for reindexing.

The steps involved in data synchronization are the following:

1. Oracle spawns the crawler according to the synchronization schedule you specify with the administration tool. The URL queue is populated with the data source URLs assigned to the schedule.
2. Crawler initiates multiple crawling threads.
3. Each crawler thread removes the next URL in the queue.

4. Each crawler thread fetches the document from the Web. The page is usually an HTML file containing text and hypertext links.
5. Each crawler thread calculates a checksum for the newly retrieved page and compares it with the checksum of the cached page. If the checksum is the same, then the page is discarded and crawler goes to step 3. Otherwise, the crawler moves to the next step.
6. Each crawler thread scans the document for hypertext links and inserts new links into the URL queue. Links that are already in the document table are discarded.
7. Crawler caches the document in the local file system. See Figure 6–2.
8. Crawler registers URL in the document table.
9. If the file system cache is full or if the URL queue is empty, then Web page caching stops and indexing begins. Otherwise, the crawler thread starts over at Step 3.

Ultra Search Remote Crawler

To increase crawling performance, set up the Ultra Search crawler to run on one or more computers separate from your database. These computers are called remote crawlers. However, each computer must share cache, log, and mail archive directories with the database computer.

To configure a remote crawler, you must first install the Ultra Search middle tier on a computer other than the database host. During installation, the remote crawler is registered with the Ultra Search system, and a profile is created for the remote crawler. After installing the Ultra Search middle tier, you must log on to the Ultra Search administration tool and edit the remote crawler profile. You can then assign a remote crawler to a crawling schedule. To edit remote crawler profiles, use the **Crawler Settings** page in the administration tool.

Note: When launching a remote crawler, the Ultra Search back end database communicates with the remote computer through Java remote method invocation (RMI). By default, RMI sends data over the network unencrypted. Using the remote crawler to perform crawling introduces a potential security risk. A malicious entity within the enterprise could steal the Ultra Search instance schema and password by listening to packets going across the network. Refrain from using the remote crawler feature if this security risk is unacceptable.

Understanding the Ultra Search Administration Tool

The Ultra Search administration tool lets you manage Ultra Search instances. This chapter helps guide you through the screens on the Ultra Search administration tool. It contains the following topics:

- Ultra Search Administration Tool
- Logging On to Ultra Search
- Logging On and Managing Instances as SSO Users
- Instances Page
- Crawler Page
- Web Access Page
- Attributes Page
- Sources Page
- Schedules Page
- Queries Page
- Users Page
- Globalization Page

Ultra Search Administration Tool

The Ultra Search administration tool is a J2EE-compliant Web application. You can use it to manage Ultra Search instances. To use the administration tool, log on as

either a database user, an Enterprise Manager super-user, a Portal user, or an SSO user through any browser.

Note: The Ultra Search administration tool and the Ultra Search query applications are part of the Ultra Search middle tier. However, the Ultra Search administration tool is independent from the Ultra Search query application. Therefore, they can be hosted on different computers to enhance security or scalability.

With the administration tool, you can do the following:

- Log on to Ultra Search
- Create Ultra Search instances
- Manage administrative users
- Define data sources and assign them to data groups
- Configure and schedule the Ultra Search crawler
- Set query options
- Translate search attributes and LOV and data group display names to different languages

Setting Crawler Parameters

To configure the Ultra Search crawler, you must do the following:

- Set crawler parameters, such as the crawler log file directory. To do so, use the Crawler Page.
- Set Web access parameters, such as authentication and the proxy server. To do so, use the Web Access Page.
- Define data sources. Data sources can be Web pages, database tables, files, email mailing lists, Oracle Sources (for example, OracleAS Portals or federated sources), or user-defined data sources. You can assign one or more data sources to a crawler schedule. To define data sources, use the Sources Page. You can also set parameters for the source, such as domain inclusions or exclusions for Web sources or the display URL template or column for table sources.

- Define synchronization schedules. The crawler uses the synchronization schedule to reconcile the Ultra Search index with current data source content. To define crawling schedules, use the Schedules Page.

Setting Query Options

Use query options to let users limit their searches. Searches can be limited to document attributes and data groups.

Attributes

Search attributes can be mapped to HTML metatags, table columns, document attributes, and email headers. Some attributes, such as author and description, are predefined and need no configuration. However, you can customize your own attributes. To set custom search attributes to expose to the query user, use the Attributes Page.

Data Groups

Data source groups are logical entities exposed to the search engine user. When entering a query, the search engine user is asked to select one or more data groups to search from. A data group consists of one or more data sources. To define data groups, use the Queries Page.

Online Help in Different Languages

Ultra Search provides context-sensitive online help, which can be viewed in different languages. You can change the language preferences in the Users Page.

Logging On to Ultra Search

The following users can log on to the Ultra Search administration tool:

- Single Sign-on (SSO) users: These users are managed by the Oracle Internet Directory (OID) and are authenticated by the SSO server. The Ultra Search administration tool identifies all Ultra Search instances to which the SSO user has access. This is available only if you have the Oracle Identity Management infrastructure installed.
- Database users (non-SSO): These users exist in the database on which Ultra Search runs.
- Enterprise Manager users

To log on to the administration tool, point your Web browser to one of the following URLs:

- For non-SSO mode:
`http://hostname:port/ultrasearch/admin/index.jsp`
- For SSO mode: `http://hostname:port/ultrasearch/admin_sso/index.jsp`

Immediately after installation, the only users able to create and manage instances are the following:

- The `WKSYS` database user
- The Enterprise Manager user
- The `PORTAL` SSO user belonging to the default company [not supported in the Oracle database release]
- The `ORCLADMIN` SSO user belonging to the default company [this is available only if the Oracle Identity Management infrastructure is installed]

After you are logged on as one of these special users, you can grant permission to other users, enabling them to create and manage Ultra Search instances. Using the Ultra Search administration tool, you can only grant and revoke Ultra Search related permissions to and from existing users. To add or delete users, use the OID for single-sign-on users or Oracle Enterprise Manager for local database users.

Note: The Ultra Search product database dictionary is installed in the `WKSYS` schema.

See Also:

- Chapter 2, "Installing and Configuring Ultra Search"
- "Changing Ultra Search Schema Passwords" on page 3-2 for information about changing the WKSYS password
- "Instances Page" on page 7-6 for more information about creating Ultra Search instances
- "Users Page" on page 7-42 for more information about granting permission to other users
- "Logging On and Managing Instances as SSO Users" on page 7-5 for more information about how Ultra Search handles SSO users

Logging On and Managing Instances as SSO Users

Note: Single Sign-On (SSO) is available only if the Oracle Identity Management infrastructure is installed

Logging On to Ultra Search

When a single sign-on (SSO) user logs on to the SSO-protected Ultra Search administration tool, the user is first prompted with the SSO login screen.

Enter the SSO user name and password. After the SSO server authenticates the user, the user sees a list of Ultra Search instances that they have the privilege to manage.

There are different URLs for different users. For example:

- SSO users: `http://<host>:<http port>/ultrasearch/admin_sso/index.jsp`
- Portal users: `http://<host>:<http port>/pls/portal`
- Enterprise Manager users: `http://<host>:<em port>/`

Granting Privileges to SSO Users

You might need to grant super-user privileges, or privileges for managing an Ultra Search instance, to an SSO user. This process is slightly different, depending on

whether OracleAS Portal is running in hosted mode or non-hosted mode, as described in the following list:

Note: An SSO user is uniquely identified by Ultra Search with an SSO-nickname/subscriber-nickname combination.

- In non-hosted mode, the subscriber-nickname is not required when granting privileges to an SSO user. This is because there is exactly one subscriber in OracleAS Portal in non-hosted mode.
- In hosted mode, the subscriber-nickname is required when granting privileges to an SSO user. This is because there can be more than one subscriber in OracleAS Portal, and two or more users with the same SSO-nickname (for example, PORTAL) could be distinct SSO users distinguished by their subscriber-nickname. When running in hosted mode, also note the following:
 - When granting permissions for the default subscriber user, always specify "DEFAULT COMPANY" for the subscriber-nickname, even though the actual nickname could be different; for example, "ORACLE". The actual nickname is not recognized by Ultra Search.
 - When logging in to SSO as the default subscriber user, leave the subscriber nickname blank. Alternatively, enter "DEFAULT COMPANY" instead of the actual subscriber nickname; for example, "ORACLE" so that it is recognized by Ultra Search.

Note: At any point after installation, you can run an OracleAS Portal script to alter the running mode from non-hosted to hosted. Whenever this is done, the OracleAS Portal script invokes an Ultra Search script to inform Ultra Search of the change from non-hosted to hosted modes.

See Also: *Hosting Developer's Guide* at <http://otn.oracle.com/>.

Instances Page

After successfully logging on to the Ultra Search administration tool, you find yourself on the **Instances Page**. This page manages all Ultra Search instances in the

local database. In the top left corner of the page, there are tabs for creating, selecting, editing, and deleting instances.

Before you can use the administration tool to configure crawling and indexing, you must create an Ultra Search instance. An Ultra Search instance is identified with a name and has its own crawling schedules and index. Only users granted super-user privileges can create Ultra Search instances.

Creating an Instance

To create an instance, click **Create**. You can create a regular instance or a read-only snapshot instance. Only users with super-user privileges can create new instances.

Note: If you define the same data source within different instances Ultra Search, then there could be crawling conflicts for table data sources with logging enabled, email data sources, and some user-defined data sources.

Creating a Regular Instance

To create an instance, do the following:

1. Prepare the database user.

Every Ultra Search instance is based on a database user/schema with the WKUSER role.

The database user you create to house the Ultra Search instance should be assigned a dedicated self-contained tablespace. This is important if you plan to ever create snapshot instances of this instance. To do this, create a new tablespace. Then, create a new database user whose default tablespace is the one you just created.

See Also:

- "Configuring the Oracle Server for Ultra Search" on page 3-2 for information and instructions on configuring database users for Ultra Search
 - "Creating a Snapshot Instance" on page 7-8
2. Follow instance creation in the Ultra Search administration tool.

From the main instance creation page, click **Create Instance**, and provide the following information:

- Instance name
- Database schema: this is the user name from step 1.
- Schema password

You can also enter the following optional index preferences:

- Lexer
Specify the name of the lexer you want to use for indexing. The lexer breaks text into tokens according to your language. These tokens are usually words. The default lexer is `wksys.wk_lexer`, as defined in the `wk0pref.sql` file. After the instance is created, the lexer can no longer be changed.
- Stoplist
Specify the name of a stoplist you want to use during indexing. The default stoplist is `wksys.wk_stoplist`, as defined in the `wk0pref.sql` file. Try to avoid modifying the stoplist after the instance has been created.
- Storage
Specify the name of the storage preference for the index of your instance. The default storage preference is `wksys.wk_storage`, as defined in the `wk0pref.sql` file. After the instance is created, the storage preference cannot be changed.

See Also:

- *Oracle Text Reference* for more information on these creating and modifying lexers, stoplists, and storage
- "Managing Stoplists" on page 3-7

Creating a Snapshot Instance

A snapshot instance is a copy of another instance. Unlike a regular instance, a snapshot instance is read only; it does not synchronize its index to the search domain. After the master instance re-synchronizes to the search domain, the snapshot instance becomes out of date. At that point, you should delete the snapshot and create a new one.

Note: The snapshot and its master instance cannot reside on the same database.

A snapshot instance is useful for the following purposes:

- Query Processing

Two Ultra Search instances can answer queries about the same search domain. Therefore, in a set amount of time, two instances can answer more queries about that domain than one instance. Because snapshot instances do not involve crawling and indexing, snapshot instance creation is fast and inexpensive. Thus, snapshot instances can improve scalability.

- Backups

If the master instance becomes corrupted, its snapshot can be transformed into a regular instance by editing the instance mode to updatable. Because the snapshot and its master instance cannot reside on the same database, a snapshot instance should be made updatable only to replace a corrupted master instance.

A snapshot instance does not inherit authentication from the master instance. Therefore, if you make a snapshot instance updatable, you must re-enter any authentication information needed to crawl the search domain.

To create a snapshot instance, do the following:

1. Prepare the database user.

As with regular instances, snapshot instances require a database user. This user must have been granted the `WKUSER` role.

2. Copy the data from the master instance.

This is done with the transportable tablespace mechanism, which does not allow renaming of tablespaces. Therefore, snapshot instances cannot be created on the same database as its master.

Identify the tablespace or the set of tablespaces that contain all the master instance data. Then, copy it, and plug it into the database user from step 1.

3. Follow snapshot instance creation in the Ultra Search administration tool.

From the main instance creation page, click **Create Read-Only Snapshot Instance**, and provide the following information:

- Snapshot instance name
 - Snapshot schema name: this is the database user from step 1.
 - Snapshot schema password
 - Database link: this is the name of the database link to the database where the master instance lives.
 - Master instance name
4. Enable the snapshot for secure searches.

If the master instance for the snapshot of is secure-search enabled and if the destination database that you are making a snapshot in supports secure-search enabled instances, then you must also run a PL/SQL procedure in the destination database where you are creating the snapshot.

Running this procedure translates the IDs of the access control lists (ACLs) in the destination database, rendering them usable. Log on to the database as the WKSYS user. Invoke the procedure as follows:

```
exec WK_ADM.USE_INSTANCE('instance_name');  
exec WK_ADM.TRANSLATE_ACL_IDS();
```

where *instance_name* is the name of the snapshot instance

Make sure that this statement completes successfully without error.

See Also:

- Chapter 3, "Post-Installation Information" for information on changing the WKSYS password and for instructions on configuring database users for Ultra Search
- *Oracle9i Database Administrator's Guide* for details on using transportable tablespaces

Selecting an Instance

You can have multiple Ultra Search instances. For example, an organization could have separate Ultra Search instances for its marketing, human resources, and development portals. The administration tool requires you to specify an instance before it lets you make any instance-specific changes.

To select an instance, do the following:

1. Click **Select** on the **Instances Page**.

2. Select an instance from the pull-down menu.
3. Click **Apply**.

Note: Instances do not share data. Data sources, schedules, and indexes are specific to each instance.

Deleting an Instance

To delete an instance, do the following:

1. Click **Delete** on the **Instances Page**.
2. Select an instance from the pull-down menu.
3. Click **Apply**.

Note: To delete an Ultra Search instance, the user must be granted the super-user privileges.

Editing an Instance

To edit an instance, click **Edit** on the **Instances Page**.

You can change the instance mode (make the instance updatable) or change the instance password.

Instance Mode

You can change the instance mode to updatable or read only. Updatable instances synchronize themselves to the search domain on a set schedule, whereas read-only instances (snapshot instances) do not do any synchronization. To set the instance mode, select the box corresponding to the mode you want, and click **Apply**.

Schema Password

An Ultra Search instance must know the password of the database user in which it resides. The instance cannot get this information directly from the database. During instance creation, Oracle provides the database user password, and the instance caches this information.

If this database user password changes, then the password that the instance has cached must be updated. To do this, enter the new password and click **Apply**. After the new password is verified against the database, it replaces the cached password.

Crawler Page

The Ultra Search crawler is a Java application that spawns threads to crawl defined data sources, such as Web sites, database tables, or email archives. Crawling occurs at regularly scheduled intervals, as defined in the Schedules Page.

With this page, you can do the following:

Configure the Settings

Crawler Threads

Specify the number of crawler threads to be spawned at run time.

Number of Processors

Specify the number of central processing units (CPUs) that exist on the server where the Ultra Search crawler will run. This setting determines the optimal number of document conversion threads used by the system. A document conversion thread converts multiformat documents into HTML documents for proper indexing.

Automatic Language Detection

Not all documents retrieved by the Ultra Search crawler specify the language. For documents with no language specification, the Ultra Search crawler attempts to automatically detect language. Click **Yes** to turn on this feature.

The language recognizer is trained statistically using trigram data from documents in various languages (Danish, Dutch, English, French, German, Italian, Portuguese, and Spanish). It starts with the hypothesis that the given document does not belong to any language and ultimately refutes this hypothesis for a particular language where possible. It operates on Latin-1 alphabet and any language with a deterministic Unicode range of characters (Chinese, Japanese, Korean, and so on).

The crawler determines the language code by checking the HTTP header content-language or the LANGUAGE column, if it is a table data source. If it cannot determine the language, then it takes the following steps:

1. If the language recognizer is not available or if it is unable to determine a language code, then the default language code is used

2. If the language recognizer is available, then the output from the recognizer is used.

This language code is populated in 'LANG' column of the `wk$url` and `wk$doc` tables. Multilexer is the only lexer used for Ultra Search. All document URLs are stored in `wk$doc` for indexing and `wk$url` for crawling.

Default Language

If automatic language detection is disabled, or if a Web document does not have a specified language, then the crawler assumes that the Web page is written in this default language. This setting is important, because language directly determines how a document is indexed.

Note: This default language is used only if the crawler cannot determine the document language during crawling. Set language preference in the Users Page.

You can select a default language for the crawler or for data sources. Default language support for indexing and querying is available for the following languages:

- Polish
- Chinese
- Hungarian
- Norwegian
- Romanian
- Finnish
- Japanese
- Spanish
- Slovak
- English
- Turkish
- Danish
- Swedish

- Russian
- German
- Korean
- Dutch
- Italian
- Greek
- Portuguese
- Czech
- Hebrew
- French
- Arabic

Crawling Depth

A Web document could contain links to other Web documents, which could contain more links. This setting lets you specify the maximum number of nested links the crawler will follow.

See Also: "Tuning the Web Crawling Process" on page 4-2 for more information on the importance of the crawling depth

Crawler Timeout Threshold

Specify in seconds a crawler timeout. The crawler timeout threshold is used to force a timeout when the crawler cannot access a Web page.

Default Character Set

Specify the default character set. The crawler uses this setting when an HTML document does not have its character set specified.

Cache Directory

Specify the absolute path of the cache directory. During crawling, documents are stored in the cache directory. Every time the preset size is reached, crawling stops and indexing starts.

If you are crawling sensitive information, then make sure that you set the appropriate file system read permission to the cache directory.

You can choose whether or not to have the cache cleared after indexing.

Crawler Logging

Specify the following:

- Level of detail: everything or only a summary
- Crawler logfile directory
- Crawler logfile language

The log file directory stores the crawler log files. The log file records all crawler activity, warnings, and error messages for a particular schedule. It includes messages logged at startup, runtime, and shutdown. Logging everything can create very large log files when crawling a large number of documents. However, in certain situations, it can be beneficial to configure the crawler to print detailed activity to each schedule log file. The crawler logfile language is the language the crawler uses to generate the log file.

Database Connect String

The database connect string is a standard JDBC connect string used by the remote crawler when it connects to the database. The connect string can be provided in the form of `[hostname]:[port]:[sid]` or in the form of a TNS keyword-value syntax; for example:

```
"(DESCRIPTION=(ADDRESS=(PROTOCOL=tcp)(HOST=...) (PORT=5521)...))"
```

See Also: *Oracle9i JDBC Developer's Guide and Reference*

You can update the JDBC connect string to a different format; for example, an LDAP format. However, you cannot change the JDBC connect string to point to a different database. The JDBC connect string must be set to the database where the middle tier points; that is, the middle tier and the JDBC should point to the same database.

In a Real Application Clusters environment, the TNS keyword-value syntax should be used, because it allows connection to any node of the system. For example,

```
"(DESCRIPTION=(LOAD_
BALANCE=yes)(ADDRESS=(PROTOCOL=TCP)(HOST=c1s02a)(PORT=3001))
(ADDRESS=(PROTOCOL=TCP)(HOST=c1s02b)(PORT=3001)))(CONNECT_DATA=(SERVICE_
NAME=sales.us.acme.com))"
```

Remote Crawler Profiles

Use this page to view and edit remote crawler profiles. A remote crawler profile consists of all parameters needed to run the Ultra Search crawler on a remote computer other than the Ultra Search database. A remote crawler profile is identified by the host name. The profile includes the cache, log, and mail directories that the remote crawler shares with the database computer.

To set these parameters, click **Edit**. Enter the shared directory paths as seen by the remote crawler. You must ensure that these directories are shared or mounted appropriately.

Crawler Statistics

Use this page to view the following crawler statistics:

Summary of Crawler Activity

This provides a general summary of crawler activity:

- Aggregate crawler statistics
- Total number of documents indexed
- Crawler statistics by data source type

Detailed Crawler Statistics

This includes the following:

- List of hosts crawled and indexed
- Document distribution by depth
- Document distribution by document type
- Document distribution by data source type

Crawler Progress

This displays crawler progress for the past week. It shows the total number of documents that have been indexed for exactly one week prior to the current time. The **Time** column rounds the current time to the nearest hour.

Problematic URLs

This lists errors encountered during the crawling process. It also lists the number of URLs that cause each error.

Web Access Page

Use this page to set up authentication and proxies.

Proxies

Specify a proxy server if the search space includes Web pages that reside outside your organization's firewall. Specifying a proxy server is optional. Currently, only the HTTP protocol is supported.

Note: The crawler cannot use a proxy server that requires proxy authentication.

You can also set domain exceptions.

Authentication

Use this page to enter authentication information global to all data sources.

Note: The data source specific authentication take precedence over this global authentication.

HTTP Authentication

Specify the user name and password for the host and realm for which HTTP authentication is required. Ultra Search supports both basic and digest authentication.

HTML Forms

Register HTML forms that you want the Ultra Search crawler to automatically fill out during Web crawling. HTML form support requires that HTTP cookie functionality is enabled.

You can register HTML forms manually or with the form registration wizard. If the HTML form contains JavaScript, then the wizard might fail and you will need to use manual registration

Attributes Page

When your indexed documents contain metadata, such as author and date information, you can let users refine their searches based on this information. For example, users can search for all documents where the author attribute has a certain value.

The list of values (LOV) for a document attribute can help specify a search query. An attribute value can have a display name for it. For example, the attribute country might use country code as the attribute value, but show the name of the country to the user. There could be multiple translations of the attribute display name.

To define a search attribute, use the **Search Attributes** subtab. Ultra Search provides some system-defined attributes, such as author and description. You can also define your own.

After defining search attributes, you must map between document attributes and global search attributes for data sources. To do so, use the **Mappings** subtab.

Note: Ultra Search provides a command-line tool to load metadata, such as search attribute LOVs and display names into an Ultra Search database. If you have a large amount of data, this is probably faster than using the HTML-based administration tool. For more information, see Appendix A, "Loading Metadata into Ultra Search".

Search Attributes

Search attributes are attributes exposed to the query user. Ultra Search provides system-defined attributes, such as author and description. Ultra Search maintains a global list of search attributes. You can add, edit, or delete search attributes. You can also click **Manage LOV** to change the list of values (LOV) for the search attribute. There are two categories of attribute LOVs: one is global across all data sources, the other is data source-specific.

To define your own attribute, enter the name of the attribute in the text box; select **string**, **date**, or **number**; and click **Add**.

You can add or delete LOV entry and display name for search attributes. Display name is optional. If display name is absent, then LOV entry is used in the query screen.

Note: LOV is only represented as string type. If LOV is in date format, then you must use "DD-MM-YYYY" to enter the LOV.

To update the policy value, click **Manage LOV** for any attribute.

A data source-specific LOV can be updated in three ways:

- Update the LOV manually.
- The crawler agent can automatically update the LOV during the crawling process.
- New LOV entries can be automatically added by inspecting attribute values of incoming documents.

Caution: If the update policy is agent-controlled, then the LOV and all translated values are erased in the next crawling.

Mappings

This section displays mapping information for all data sources. For user-defined data sources, mapping is done at the agent level, and document attributes are automatically mapped to search attributes with the same name initially. Document attributes and search attributes are mapped one-to-one. For each user-defined data source, you can edit the global search attribute to which the document attribute is mapped.

For Web, file, or table data sources, mappings are created manually when you create the data source. For user-defined data sources, mappings are automatically created on subsequent crawls.

Click **Edit Mappings** to change this mapping.

Editing the existing mapping is costly, because the crawler must recrawl all documents for this data source. Avoid this step, unless necessary.

Note: There are no user-managed mappings for email sources. There are two predefined mappings for emails. The "From" field of an email is intrinsically mapped to the Ultra Search author attribute. Likewise, the "Subject" field of an email is mapped to the Ultra Search subject attribute. The abstract of the email message is mapped to the description attribute.

Sources Page

A collection of documents is called a source. The data source is characterized by the properties of its location, such as a Web site or an email inbox. The Ultra Search crawler retrieves data from one or more data sources.

The different types of sources are:

- Web Sources
- Table Sources
- Email Sources
- File Sources
- Oracle Sources
- User-Defined Sources (requires a crawler agent)

See Also:

- "Schedules Page" on page 7-32 to assign one or more data sources to a synchronization schedule
- "Queries Page" on page 7-38 to assign data sources to data groups to enable restrictive querying

You can create as many data sources as you want. The following section explains how to create and edit data sources.

Web Sources

A Web source represents the content on a specific Web site. Web sources facilitate maintenance crawling of specific Web sites.

Creating Web Sources

To create a new Web source, do the following:

1. Specify a name for the Web source and a starting address. This is the URL for the crawler to begin crawling. The starting address can be HTTP or HTTPS.
2. Set URL boundary rules to refine the crawling space. You can include or exclude hosts or domains beginning with, ending with, or equal to a specific name.

For example, an inclusion domain ending with `oracle.com` limits the Ultra Search crawler to hosts belonging to Oracle worldwide. Anything ending with `oracle.com` is crawled; but, `http://www.oracle.com.tw` is not crawled. If you change the inclusion domain to `yahoo.com` with a new seed "`http://www.yahoo.com`", then all `oracle.com` URLs are dropped by the crawler.

An exclusion domain `uk.oracle.com` prevents the crawler from crawling Oracle hosts in the United Kingdom. You can also include or exclude Web sites with a specific port. (By default, all ports are crawled.) You can have port inclusion or port exclusion rules for a specific host, but not both. Exclusion rules always override inclusion rules.

3. Specify the types of documents the Ultra Search crawler should process for this source. HTML and plain text are default document types that the crawler always processes.
4. Specify the authentication settings. This step is optional. Under **HTTP Authentication**, specify the user name and password for host-realm for which authentication is required. Under **HTML Forms**, you can register HTML forms that you want the Ultra Search crawler to automatically fill out during Web crawling. HTML form support requires that HTTP cookie functionality is enabled. Click **Register HTML Form** to register authentication forms protecting the data source. Note: For the form URL to be crawled, you must verify that the URL is not excluded in the `robots.txt` file. If so, then you must disable robot exclusion for this data source. (By default, Ultra Search enables robot exclusion.)
5. Choose either **No ACL** or **Ultra Search ACL** for the data source. When a user performs a search, the ACL (access control list) controls which documents the user can access. The default is no ACL, with all documents considered searchable and visible. You can add more than one group and user to the ACL for the data source. The option to choose is only available if the instance is security-enabled.

6. Define, edit, or delete metatag mappings for your Web source. Metatags are descriptive tags in the HTML document header. One metatag can map to only one search attribute.
7. Override the default crawler settings for each Web source. This step is optional. The parameters you can override are the crawling depth, the number of crawler threads, the language, the crawler timeout threshold, the character set, the maximum cookie size, the maximum number of cookies, and the maximum number of cookies for each host. You can also enable or disable robots exclusion, language detection, the `UrlRewriter`, indexing dynamic pages, HTTP cookies, and whether content of the cookie log file is shown. (You can also edit those in **Edit Web Sources**.)

Robots exclusion lets you control which parts of your sites can be visited by robots. If robots exclusion is enabled (default), then the Web crawler traverses the pages based on the access policy specified in the Web server `robots.txt` file. For example, when a robot visits `http://www.foobar.com/`, it checks for `http://www.foobar.com/robots.txt`. If it finds it, the crawler analyzes its contents to see if it is allowed to retrieve the document. If you own the Web sites, then you can disable robots exclusions. However, when crawling other Web sites, you should always comply with `robots.txt` by enabling robots exclusion.

The URL Rewriter is a user-supplied Java module for implementing the Ultra Search `UrlRewriter` interface. It is used by the crawler to filter or rewrite extracted URL links before they are put into the URL queue. URL filtering removes unwanted links, and URL rewriting transforms the URL link. This transformation is necessary when access URLs are used.

The `UrlRewriter` provides the following possible outcomes for links:

- There is no change to the link. The crawler inserts it as it is.
- Discard the link. There is no insertion.
- A new display URL is returned, replacing the URL link for insertion.
- A display URL and an access URL are returned. The display URL may or may not be identical to the URL link.

The generated new URL link is subject to all existing host, path, and mimetype inclusion and exclusion rules.

You must put the implemented rewriter class in a jar file and provide the class name and jar file name here.

If Index Dynamic Page is set to **Yes**, then dynamic URLs are crawled and indexed. For data sources already crawled with this option, setting Index Dynamic Page to **No** and recrawling the data source removes all dynamic URLs from the index.

Some dynamic pages appear as multiple search hits for the same page, and you may not want them all indexed. Other dynamic pages are each different and need to be indexed. You must distinguish between these two kinds of dynamic pages. In general, dynamic pages that only change in menu expansion without affecting its contents should not be indexed. Consider the following three URLs:

```
http://itweb.oraclecorp.com/aboutit/network/npe/standards/naming_
convention.html
```

```
http://itweb.oraclecorp.com/aboutit/network/npe/standards/naming_
convention.html?nsdnv=14z1
```

```
http://itweb.oraclecorp.com/aboutit/network/npe/standards/naming_
convention.html?nsdnv=14
```

The question mark (?) in the URL indicates that the rest of the strings are input parameters. The duplicate hits are essentially the same page with different side menu expansion. Ideally, the same query should yield only one hit:

```
http://itweb.oraclecorp.com/aboutit/network/npe/standards/naming_
convention.html
```

Dynamic page index control applies to the whole data source. So, if a Web site has both kinds of dynamic pages, you need to define them separately as two data sources in order to control the indexing of those dynamic pages.

See Also:

- "Ultra Search URL Rewriter API" on page 8-29
- "Using Crawler Agents" on page 6-3
- "Crawler Page" on page 7-12 for information on default languages

Table Sources

A table source represents content in a database table or view. The database table or view can reside in the Ultra Search database instance or in a remote database. Ultra Search accesses remote databases using database links.

See Also: "Limitations With Database Links" on page 7-25

Creating Table Sources

To create a table source, click **Create Table Source**, and follow these steps:

1. Specify a table source name, and the name of the database link, schema, and table. Click **Locate Table**.
2. Specify settings for your table source, such as the default language and the primary key column. You can also specify the column where final content should be delivered, and the type of data stored in that column; for example, HTML, plain text, or binary. For information on default languages, see "Crawler Page" on page 7-12.
3. Verify the information about your table source.
4. Decide whether or not to use the Ultra Search logging mechanism to optimize the crawling of table data sources. When crawling is enabled, only newly updated documents are revisited during the crawling process. You can enable logging for Oracle tables, enable logging for non-Oracle tables, or disable the logging mechanism. If you enable logging, then you are prompted to create a log table and log triggers. Oracle SQL statements are provided for Oracle tables. If you are using non-Oracle tables, then you must manually create a log table and log triggers. Follow the examples provided to create the log table and log triggers. After you have created the table, enter the table name in **Log Table Name**.
5. Map table columns to search attributes. Each table column can be mapped to exactly one search attribute. This lets the search engine seamlessly search data from the table source.
6. Specify the display URL template or column for the table source. This step is optional. Ultra Search uses a default text viewer for table data sources. If you specify display URL, then Ultra Search uses the Web URL defined to display the table data retrieved. If display URL column is available, then Ultra Search uses the column to get the URL to display the table data source content. You can also specify display URL templates in the following format:
`http://hostname:port/path?parameter_name=$(key1)` where `key1` is the corresponding table's primary key column. For example, assume that you can use the following URL to query the bug number 1234567, and the bug number is the primary key of the table:
`http://bug:7777/pls/bug?rptno=1234567`. You can set the table source display URL template to `http://bug:7777/pls/bug?rptno=$(key1)`.

The **Table Column to Key Mappings** section provides mapping information. Ultra Search supports table keys in `STRING`, `NUMBER`, or `DATE` type. If key1 is of `NUMBER` or `DATE` type, then you must specify the format model used by the Web site so that Oracle knows how to interpret the string. For example, the date format model for the string '11-Nov-1999' is 'DD-Mon-YYYY'. You can also map other table columns to Ultra Search attributes. Do not map the text column.

7. Specify the ACL (access control list) policy for the data source. When a user performs a search, the ACL controls which documents the user can access. The default is no ACL, with all documents considered public and visible. Alternatively, you can specify to use Ultra Search ACL. You can add more than one group and user to the ACL for the data source. The option to choose is only available if the instance is security-enabled.

See Also: *Oracle9i SQL Reference* for more on format models

Editing Table Sources

On the main **Table Sources** page, click **Edit** to change the name of the table source. You can change, add, or delete table column and search attribute mappings; change the display URL template or column; and view values of the table source settings.

Table Sources Comprised of More Than One Table

If a table source has more than one table, then a view joining the relevant tables must be created. Ultra Search then uses this view as the table source. For example, two tables with a master-detail relationship can be joined through a `SELECT` statement on the master table and a user-implemented PL/SQL function that concatenate the detail table rows.

Limitations With Database Links

The following restrictions apply to base tables or views on a remote database that are accessed over a database link by the crawler.

- If the text column of the base table or view is of type `BLOB` or `CLOB`, then the table must have a `ROWID` column. A table or view might not have a `ROWID` column for various reasons, including the following:
 - A view is comprised of a join of one or more tables.
 - A view is based on a single table using a `GROUP BY` clause.

The best way to know if a remote table or view can be safely crawled by Ultra Search is to check for the existence of the `ROWID` column. To do so, run the following SQL statement against that table or view using SQL*Plus:

```
SELECT MIN(ROWID) FROM table_name/view_name;
```

- The base table or view cannot have text columns of type `BFILE`, `RAW`.

Email Sources

An email source derives its content from emails sent to a specific email address. When the Ultra Search crawler searches an email source, it collects all emails that have the specific email address in any of the "To:" or "Cc:" email header fields.

The most popular application of an email source is where an email source represents all emails sent to a mailing list. In such a scenario, multiple email sources are defined where each email source represents an email list.

To crawl email sources, you need an IMAP account. At present, the Ultra Search crawler can only crawl one IMAP account. Therefore, all emails to be crawled must be found in the inbox of that IMAP account. For example, in the case of mailing lists, the IMAP account should be subscribed to all desired mailing lists. All new postings to the mailing lists are sent to the IMAP email account and subsequently crawled. The Ultra Search crawler is IMAP4 compliant.

When the Ultra Search crawler retrieves an email message, it deletes the email message from the IMAP server. Then, it converts the email message content to HTML and temporarily stores that HTML in the cache directory for indexing. Next, the Ultra Search crawler stores all retrieved messages in a directory known as the archive directory. The email files stored in this directory are displayed to the search end-user when referenced by a query hit.

To crawl email sources, you must specify the user name and password of the email account on the IMAP server. Also specify the IMAP server host name and the archive directory.

Creating Email Sources

To create email sources, you must enter an email address and a description. Optionally, you can specify email aliases and ACL policy. The description can be viewed by all search end-users, so you should specify a short but meaningful name. When you create (register) an email source, the name you use is the email of the mailing list. If the emails are not sent to one of the registered mailing lists, then those emails are not crawled.

You can specify email address aliases for an email source. Specifying an alias for an email source causes all emails sent to the main email address, as well as the alias address, to be gathered by the crawler. An alias is useful when two or more email addresses are logically the same. For example, an email source representing the distribution list `list@company.com` can have the alternate address `list@my.company.com`. If `list@my.company.com` is added to the alias list, then email sent to that address are treated as if they were sent to `list@company.com`.

Specify the ACL (access control list) policy for the data source. When a user performs a search, the ACL controls which documents the user can access. The default is no ACL, with all documents considered searchable and visible. You can add more than one group and user to the ACL for the data source.

File Sources

A file source is the set of documents that can be accessed through the file protocol on the local machine.

To edit the name of a file source, click **Edit**.

Creating File Sources

To create a new file source, do the following:

1. Specify a name for the file source and the default language.
2. Designate files or directories to be crawled. If a URL represents a single file, then the Ultra Search crawler searches only that file. If a URL represents a directory, then the crawler recursively crawls all files and subdirectories in that directory.
3. Specify inclusion and exclusion paths to modify the crawling space associated with this file source. This step is optional. An inclusion path limits the crawling space. An exclusion path lets you further define the crawling space. If neither path is specified, then crawling is limited to the underlying file system access privileges. Path rules are host-specific, but you can specify more than one path rule for each host. For example, on the same host, you can include path `files://host/doc` and exclude path `files://host/doc/unwanted`.
4. Specify the types of documents the Ultra Search crawler should process for this file source. HTML and plain text are default document types that the crawler always processes.

5. Ultra Search displays file data sources in text format. However, if you specify display URL for the file data source, then Ultra Search uses the URL to display the file data source.

With display URL for file data sources, the URL uses network protocols, such as HTTP or HTTPS, to access the file data source. To generate display URL for the file data source, specify the prefix of the original file URL and the prefix of the display URL. Ultra Search replaces the prefix of the file URL with the prefix of the display URL.

For example, if your file URL is `file:///home/operation/doc/file.doc` and the display URL is `https://webhost/client/doc/file.doc`, then you can specify the file URL prefix to `file:///home/operation` and the display URL prefix to `https://webhost/client`.

6. Specify the ACL (access control list) policy for the data source. When a user performs a search, the ACL controls which documents the user can access. The default is no ACL, with all documents considered searchable and visible. Alternatively, you can specify to use Ultra Search ACL. You can add more than one group and user to the ACL for the data source. The option to choose is only available if the instance is security-enabled.

Oracle Sources

You can create, edit, or delete Oracle sources. You can choose federated or OracleAS Portal (crawlable) data sources. A federated source is a repository that maintains its own index. Ultra Search can issue a query, and the repository can return query results. Ultra Search also supports the crawling and indexing of OracleAS Portal installations. This enables searching across multiple portal installations.

Oracle Portal Sources

To create Portal sources, you must first register your portal with Ultra Search. To register your portal:

1. Provide a name and portal URL base. The portal name is used to identify this portal entry in the **Oracle Portal List** page. The URL base is the beginning portion of the portal homepage. This include host name, port number, and DAD. After it is created, the portal URL base is not updatable. Click **Register Portal**. Ultra Search attempts to contact the OracleAS Portal instance and retrieve information about it.
2. Choose one or more page groups for indexing. A portal data source is created for each page group. Click **Delete** to remove existing portal data sources.

You can edit the types of documents the Ultra Search crawler should process for a portal source. HTML and plain text are default document types that the crawler always processes. To edit document types, click **Edit** for the portal source after it has been created.

See Also: The OracleAS Portal documentation.

Federated Sources

To create federated sources, specify the name and JNDI for the new data source. By default, no resource adapter is available.

To create a federated source, you must manually deploy the Ultra Search resource adapter, or *searchlet*. A searchlet is a Java module deployed in the middle tier (inside OC4J) that searches the data in an enterprise information system on behalf of a user. A searchlet is a Java module deployed in the middle tier (inside OC4J) that searches the data in an enterprise information system on behalf of a user. When a user's query is delegated to the searchlet, the searchlet runs the query on behalf of the user. Every searchlet is a JCA 1.0 compliant resource adapter.

See Also: The JCA 1.0 spec from Javasoft for detailed information on resource adapters and Java Connector Architecture

Deploying and Binding the Ultra Search Searchlet The Ultra Search searchlet enables queries against one Ultra Search instance. The Ultra Search searchlet is packaged as `ultrasearch.rar` and is shipped under the `$ORACLE_HOME/ultrasearch/adapter/` directory.

To deploy the Ultra Search searchlet in OC4J standalone, use `admin.jar` as follows:

```
java -jar admin.jar ormi://<hostname> <admin> <welcome> -deployconnector -file ultrasearch.rar -name UltraSearchSearchlet
```

At this point, `ultrasearch.rar` has been deployed in OC4J. However, it has not been instantiated to connect to any Ultra Search instance. The Ultra Search searchlet can be instantiated multiple times, to connect to several Ultra Search instances, by repeating the following steps.

To instantiate the searchlet, configuration parameters values must be specified, and a JNDI location must be specified where the searchlet instance should be bound to. To do this, you must manually edit `oc4j-ra.xml`. This file is typically located under the `$J2EE_`

HOME/application-deployments/default/UltraSearchSearchlet/directory. The Ultra Search searchlet requires four configuration properties: connectionURL, userName, password, and instanceName. For example, to bind a searchlet under "eis/UltraSearch" to connect to the default instance 'wk_test' on machine 'dbhost', the following entry can be used:

```
<connector-factory location="eis/UltraSearch" connector-name="Ultra Search Adapter">
  <config-property name="connectionURL"
value="jdbc:oracle:thin:@dbhost:1521:sid"/>
  <config-property name="userName" value="wk_test"/>
  <config-property name="password" value="wk_test"/>
  <config-property name="instanceName" value="wk_test"/>
</connector-factory>
```

After editing oc4j-ra.xml, restart the OC4J instance. If you do not see errors upon restart, then the searchlet has been successfully instantiated and bound to JNDI.

Deploying and Binding the Federator Searchlet The Federator searchlet interacts with other searchlets to provide a single point of search against multiple repositories. For example, the Federator searchlet can invoke multiple Ultra Search searchlets to simultaneously query against multiple Ultra Search instances. In the same manner, the Federator searchlet can invoke searchlets for Oracle Files, Email, and so on.

The Federator searchlet is configured and managed with the Ultra Search administration tool, under the **Federated Sources** tab.

The Federator searchlet is packaged as federator.rar and is shipped under the \$ORACLE_HOME/ultrasearch/adapter/ directory.

The deployment procedure for federator.rar is similar to the deployment of the Ultra Search searchlet. To deploy the Federator searchlet in OC4J standalone, use admin.jar as follows:

```
java -jar admin.jar ormi://<hostname> <admin> <welcome> -deployment -file
federator.rar -name FederatorSearchlet
```

To instantiate the searchlet, the Federator searchlet requires four configuration properties: connectionURL, userName, password, and instanceName in the oc4j-ra.xml file. This file is typically located under the \$J2EE_HOME/application-deployments/default/FederatorSearchlet/directory. For example:

```
<connector-factory location="eis/Federator" connector-name="Federator Adapter">
```

```
<config-property name="connectionURL"
value="jdbc:oracle:thin:@dbhost:1521:sid"/>
<config-property name="userName" value="wk_test"/>
<config-property name="password" value="wk_test"/>
<config-property name="InstanceName" value="wk_test"/>
</connector-factory>
```

After editing `oc4j-ra.xml`, restart the OC4J instance. If you do not see errors upon restart, then the searchlet has been successfully instantiated and bound to JNDI.

User-Defined Sources

Ultra Search lets you define, edit, or delete your own data sources and types in addition to the ones provided. You might implement your own crawler agent to crawl and index a proprietary document repository or management system, such as Lotus Notes or Documentum, which contain their own databases and interfaces.

For each new data source type, you must implement a crawler agent as a Java class. The agent collects document URLs and associated metadata from the proprietary document source and returns the information to the Ultra Search crawler, which enqueues it for later crawling.

See Also: "Ultra Search Crawler Agent API" on page 8-19

To define a new data source, you first define a data source type to represent it.

Creating User-Defined Data Source Types

To create, edit, or delete data source types, click **Manage Source Types**. To create a new type, click **Create New Type**.

1. Specify data source type name, description, and crawler agent Java class file or jar file name. The crawler agent Java classpath is predefined at installation time. The agent collects the list of document URLs and associated metadata from the proprietary document source and returns it to the Ultra Search crawler, which enqueues the information for later crawling. The agent class file or jar file must be located under `$ORACLE_HOME/ultrasearch/lib/agent/`.
2. Specify parameters for this data source type. If you add parameters, you must enter the parameter name and a description. Also, you must decide whether to encrypt the parameter value.

Edit data source type information by changing the data source type name, description, crawler agent Java class/jar file name, or parameters.

Creating User-Defined Sources

To create a user-defined data source, select the type and click **Go**

1. Specify a name, default language, and parameter values for the data source. For information on default languages, see the [Crawler Page](#).
2. Specify the authentication settings. This step is optional. Under **HTTP Authentication**, specify the user name and password for host and realm for which authentication is required. Under **HTML Forms**, you can register HTML forms that you want the Ultra Search crawler to automatically fill out during Web crawling. HTML form support requires that HTTP cookie functionality is enabled. Click **Register HTML Form** to register authentication forms protecting the data source.
3. Specify the ACL (access control list) policy for the data source: no ACL, repository-generated ACL, or Ultra Search ACL. When a user performs a search, the ACL controls which documents the user can access. The default is no ACL, with all documents considered searchable and visible. For the Ultra Search ACL, you can add more than one group and user to the ACL for the data source.
4. Specify mappings. This step is optional. Document attributes are automatically mapped directly to the search attribute with the same name during crawling. If you want document attributes to map to another search attribute, then you can specify it here. The crawler picks up attributes that have been returned by the crawler agent or specified here.
5. Edit crawling parameters.
6. Specify the document types that the crawler should process for this data source. By default, HTML and plain text are always processed.

You can edit user-defined data sources by changing the name, type, default language, or starting address.

Schedules Page

Use this page to schedule data synchronization and index optimization. Data synchronization means keeping the Ultra Search index up to date with all data sources. Index optimization means keeping the updated index optimized for best query performance.

See Also: "Synchronizing Data Sources" on page 6-3

Data Synchronization

The tables on this page display information about synchronization schedules. A synchronization schedule has one or more data sources assigned to it. The synchronization schedule frequency specifies when the assigned data sources should be synchronized. Schedules are sorted first by name. Within a synchronization schedule, individual data sources are listed and can be sorted by source name or source type.

Creating Synchronization Schedules

To create a new schedule, click **Create New Schedule** and follow these steps:

1. Name the schedule.
2. Pick a schedule frequency and determine whether the schedule should automatically accept all URLs for indexing or examine URLs before indexing. You can also associate the schedule with a remote crawler profile.
3. Assign data sources to the schedule. After a data source has been assigned to a group, it cannot be assigned to other groups.

Updating Schedules

Update the indexing option in the **Update Schedule** page.

Editing Synchronization Schedules

After a synchronization schedule has been defined, you can do the following in the **Synchronization Schedules List**:

- To assign the schedule to either a crawler that runs on the database host or a remote crawler that runs on a separate host, click **Hostname**.
- To change its frequency, click the schedule interval text.
- To alter its status, click **Status**.
- To delete it, click **Delete**.
- To edit its name, data source assignments, recrawl policy, or crawling mode, click **Edit**. When the crawler retrieves a document, it checks to see if it has changed. By default, if the document has not changed, the crawler does not

process it. In certain situations, you might want to force the crawler to reprocess all documents. Click **Edit** to edit schedules in the following ways:

- Update schedule name. This step is optional. To change the schedule name, specify a name for the schedule, and click **Update Schedule Name**.
- Assign data sources to schedule. To assign a data source, select one or more available sources and click >>. After a data source has been assigned to a group, it cannot be assigned to any other group. To undo assignments of a data source, select one or more scheduled sources and click <<.
- Update crawler recrawl policy. You can update the recrawl policy to the following:
 - * **Process Documents That Have Changed:** This is maintenance crawling. Only documents that have changed are recrawled and indexed. For Web data sources, if there are new links in the updated document, then they are followed. For file data sources, new files are collected if its parent directory has changed.
 - * **Process All Documents:** The crawler recrawls the data source. For example, suppose you want to crawl only text and HTML on a Web site. Later, you also want to crawl Microsoft Word and Adobe PDF documents. You must modify the document types for the source, edit the schedule to select **Process All Documents**, then reexecute the schedule so that the crawler picks up PDF and doc document types for this data source. The crawler treats every document as if it has been changed, which means each document is fetched and processed again.

Upon relaunching the schedule, the following rules determine which URLs will be recrawled:

- * If the previous crawl did not finish (for example, you stopped the crawling or the database tablespace was full), then the crawler only crawls URLs left in the URL queue. URLs already crawled are not touched on recrawl.
- * If the URL queue is empty but there is a new seed added since the last crawl, then the crawler only crawls the new seed.
- * If the URL queue is empty and there is no new seed URL, then the crawler recrawls all crawled URLs.

Therefore, if you stop the crawler and set Index Dynamic Pages to **No**, this only affects the URLs in the queue yet to be crawled. The already crawled

dynamic pages are removed from the index on the third recrawl when the queue is empty.

Note: All crawled URLs are subject to crawler setting enforcement, not just newly crawled URLs.

- Update crawling mode. You can update the crawling mode to the following:
 - * Automatically accept all URLs for indexing: This mode crawls and indexes.
 - * Examine URLs before indexing: This mode crawls only.
 - * Index only: This mode indexes only.

The crawler behaves differently for the documents collected.

Crawling mode and recrawl policy can be combined for six different combinations. For example, **Process All Documents** and **Index Only** forces reindexing existing documents in this data source, while **Process Documents That Have Changed** and **Index Only** re-indexes only changed documents.

Launching Synchronization Schedules

A schedule's synchronization frequency can be identical to another schedule's synchronization frequency. This gives you maximum flexibility in managing data source synchronization.

You can launch a synchronization schedule in the following ways:

- Set a schedule frequency and wait for the predetermined launch time.
- Execute it immediately. To do so, click **Status**, then **Execute Immediately**.
- Manually start the schedule.

Note: Launching a synchronization schedule could take a very long time. If a schedule has been launched before, then the next time a schedule is launched, all URLs that belong to the data source to be crawled by the schedule are updated to put into a queue. Depending on the number of URLs associated with that data source, the enqueue operation may take a long time. The administration tool displays the schedule state as 'Launching' the entire time.

The launch of a schedule does not perform any enqueue if the URL queue is not empty or if there is a new seed added since the last crawl. For example, if the user stopped the crawler earlier or if the crawler terminated because of insufficient Oracle table space, then the URL queue is not empty. So, on the next launch the crawler does not try to enqueue; instead it works on the existing URL queue until it is empty. In other words, enqueue is only performed when the queue is empty at launch time.

Synchronization Status and Crawler Progress

Click the link in the status column to see the synchronization schedule status. To see the crawling progress for any data source associated with this schedule, click **Statistics**.

If you decide to examine URLs before indexing for the schedule, then after you run the schedule, the schedule status is shown as "Indexing Pending".

In data harvesting mode, you should begin crawling first. After crawling is done, click **Examine URL** to examine document URLs and status, remove unwanted documents, and start indexing. After you click **Begin Index**, you see schedule status change from launching, executing, scheduled, and so on.

The crawling progress contains the following information:

- Data source type
- Data source name
- Start time
- Finish time
- Elapsed time
- Total indexing time

- Total size of document data collected
- Average document size
- Average fetch throughput

It also contains the following statistics:

- Documents to fetch
- Documents fetched: This is the sum of Document non-indexable, Document conversion failure, and Documents indexed.
- Document fetch failures: This could be an Oracle HTTP Server timeout or another HTTP server error.
- Documents rejected: The document is not within the URL boundary rule.
- Documents discovered: This is the sum of Documents to fetch, Documents fetched, Document fetch failures, and Documents rejected.
- Documents indexed
- Documents non-indexable: This could be a file directory, a portal page that is a discovery node, or a robot metatag that specifies no index.
- Document conversion failures: The binary file filter failed.

Index Optimization

Index Optimization

To ensure fast query results, the Ultra Search crawler maintains an active index of all documents crawled over all data sources. This lets you schedule when you would like the index to be optimized. The index should be optimized during hours of low usage.

Note: Increasing the crawler temporary directory size can reduce index fragmentation.

Index Optimization Schedule

You can specify the index optimization schedule frequency. Be sure to specify all required data for the option that you select. You can optimize the index immediately, or you can enable the schedule.

Optimization Process Duration

Specify a maximum duration for the index optimization process. The actual time taken for optimization does not exceed this limit, but it could be shorter. Specifying a longer optimization time results in a more optimized index. Alternatively, you can specify that the optimization continue until it is finished.

If your Ultra Search instance is secure-search enabled, then the index optimization process also triggers garbage collection of unused access control lists (ACLs).

Queries Page

This section lets you specify query-related settings, such as data source groups, URL submission, relevancy boosting, and query statistics.

Data Groups

Data groups are logical entities exposed to the search engine user. When entering a query, the user is asked to select one or more data groups from which to search.

A data group consists of one or more data sources. A data source can be assigned to multiple data groups. Data groups are sorted first by name. Within each data group, individual data sources are listed and can be sorted by source name or source type.

To create a new data source group, do the following:

1. Specify a name for the group.
2. Assign data sources to the group. To assign a Web or table data source to this data group, select one or more available Web sources or table sources and click >>. After a data source has been assigned to a group, it cannot be assigned to any other group. To unassign a Web or table data source, select one or more scheduled sources and click <<.
3. Click **Finish**.

URL Submission

URL Submission Methods

URL submission lets query users submit URLs. These URLs are added to the seed URL list and included in the Ultra Search crawler search space. You can allow or disallow query users to submit URLs.

URL Boundary Rules Checking

URLs are submitted to a specific Web data source. URL boundary rules checking ensures that submitted URLs comply with the URL boundary rules of the Web data source. You can allow or disallow URL boundary rules checking.

Relevancy Boosting

Relevancy boosting lets administrators override the search results and influence the order that documents are ranked in the query result list. This can be used to promote important documents to higher scores. It also makes them easier to find.

See Also: "Document Relevancy Boosting" on page 1-10

There are two methods for locating URLs for relevancy boosting: locate by search or manual URL entry.

Locate by Search

To boost a URL, first locate a URL by performing a search. You can specify a host name to narrow the search. After you have located the URL, click **Information** to edit the query string and score for the document.

Manual URL Entry

If a document has not been crawled or indexed, then it cannot be found in a search. However, you can provide a URL and enter the relevancy boosting information with it. To do so, click **Create**, and enter the following:

1. Specify the document URL. You must assign the URL to a data source. This document is indexed the next time it is crawled.
2. Enter scores in the range of 1 to 100 for one or more query strings. When a user performs a search using the exact query string, the score applies for this URL.

The document is searchable after the document is loaded for the term. The document is also indexed the next time the schedule is run.

With manual URL entry, you can only assign URLs for Web data sources. Users get an error message on this page if no Web data source is defined.

Note: Ultra Search provides a command-line tool to load metadata, such as document relevance boosting, into an Ultra Search database. If you have a large amount of data, this is probably faster than using the HTML-based administration tool. For more information, see Appendix A, "Loading Metadata into Ultra Search".

Query Statistics

Enabling Query Statistics

This section lets you enable or disable the collection of query statistics. The logging of query statistics reduces query performance. Therefore, Oracle recommends that you disable the collection of query statistics during regular operation.

Note: After you enable query statistics, the table that stores statistics data is truncated every Sunday at 1:00 A.M.

Viewing Statistics

If query statistics is enabled, you can click one of the following categories:

- Daily Summary of Query Statistics
- Top 50 Queries
- Top 50 Ineffective Queries
- Top 50 Failed Queries

Daily Summary of Query Statistics This summarizes all query activity on a daily basis. The statistics gathered are:

- Average query time: the average time taken over all queries
- Number of queries: the total number of queries made in the day
- Number of hits: the average number of results returned by each query

Top 50 Queries This summarizes the 50 most frequent queries in the past 24 hours.

- Query string: the query string

- Average query time: the average time to return a result
- Number of queries: the total number of queries in the past 24 hours
- Number of hits: the average number of results returned by each query
- Frequency: the number of queries divided by total number of queries over all query strings
- Percentage of ineffective queries: the number of ineffective queries divided by total number of queries over all query strings

Top 50 Ineffective Queries This summarizes the 50 most frequent queries in the past 24 hours. Each row in the table describes statistics for a particular query string.

- Query string: the query string
- Number of queries: the total number of queries made in the past 24 hours
- Percentage of ineffective queries: the number of ineffective queries divided by total number of queries for that string

Top 50 Failed Queries This summarizes the top 50 queries that failed over the past 24 hours. A failed query is one where the search engine end-user did not locate any query results.

The columns are:

- Query string: the query string
- Number of queries: the total number of queries made in the past 24 hours
- Frequency: the percentage occurrence of a failed query
- Cumulative frequency: the cumulative percentage occurrence of all failed queries

See Also: "Tuning Query Performance" on page 4-3

Configuration

You can configure the query application and the federation engine with several parameters, including the maximum number of hits and enabling relevancy boosting.

Users Page

Use this page to manage Ultra Search administrative users. You can assign a user to manage an Ultra Search instance. You can also select a language preference.

Preferences

This section lets you set preference options for the Ultra Search administrator.

You can specify the date and time format. The pull-down menu lists the following languages:

- English
- Brazilian Portuguese
- French
- German
- Italian
- Japanese
- Korean
- Simplified Chinese
- Spanish
- Traditional Chinese

You can also select the number of rows to display on each page.

Super-Users

A user with super-user privileges can perform all administrative functions on all instances, including creating instances, dropping instances, and granting privileges. Only super-users can access this page.

Single sign-on (SSO) users can use a delegated administrative service (DAS) list of values to add another SSO user as a super-user. These users are authenticated by the SSO server before allowing access. Database users can add another database user as a super-user.

To grant super-user administrative privileges to another user, enter the user name of the user. Specify also whether the user should be allowed to grant super-user privileges to other users. Then click **Add**.

Privileges

Only instance owners, users that have been granted general administrative privileges on this instance, or super-users are allowed to access this page. Instance owners must have been granted the `WKUSER` role.

Single sign-on (SSO) users can use a delegated administrative service (DAS) list of values to add privileges to another SSO user. These users are authenticated by the SSO server before allowing access. Database users can add privileges to another database user.

Note: Database users cannot grant privileges to SSO users, and SSO users cannot grant privileges to database users. The DAS list of values only shows SSO users.

Granting general administrative privileges to a user allows that user to modify general settings for this instance. To do this, enter the user name and specify whether the user should be allowed to grant administrative privileges to other users. Then click **Add**.

To remove one or more users from the list of administrators for this instance, select one or more user names from the list of current administrators and click **Remove**.

Note: General administrative privileges do not include the ability to create or delete an instance. These privileges belong to super-users.

See Also: "Step 4: Create and Configure New Database Users for Each Ultra Search Instance" on page 3-5

Globalization Page

Ultra Search lets you translate names to different languages. This page lets you enter multiple values for search attributes, list of values (LOV) display names, and data groups.

Search Attribute Name

This section lets you translate attribute display names to different languages. The pull-down menu lists the following languages:

- English
- Arabic
- Brazilian Portuguese
- Canadian French
- Czech
- Danish
- Dutch
- Finnish
- French
- German
- Greek
- Hebrew
- Hungarian
- Italian
- Japanese
- Korean
- Latin American Spanish
- Norwegian
- Polish
- Portuguese
- Romanian
- Russian
- Simplified Chinese
- Slovak
- Spanish
- Swedish
- Thai
- Traditional Chinese

- Turkish

LOV Display Name

This section lets you translate data group names to different languages. Select a search attribute from the pull-down menu: author, description, mimetype, subject, or title. Select the LOV type, and then select the language from the pull-down menu.

Data Group Name

This section lets you translate data group display names to different languages. The pull-down menu lists the language options.

Ultra Search Developer's Guide and API Reference

This chapter explains the Ultra Search APIs and related information. This chapter contains the following topics:

- Overview of Ultra Search APIs
- Ultra Search Query API
- Customizing the Query Syntax Expansion
- Ultra Search Query Tag Library
- Ultra Search Crawler Agent API
- Ultra Search Java Email API
- Ultra Search URL Rewriter API
- Ultra Search Sample Query Applications

See Also: *Oracle Ultra Search API Reference*

Overview of Ultra Search APIs

Ultra Search provides the following APIs:

- The query API works with indexed data. The Java API does not impose any HTML rendering elements. The application can completely customize the HTML interface.
- The crawler agent API crawls and indexes proprietary document repositories.
- The email API is used by the Ultra Search query application to display emails. It can also be used when building your own custom query application.
- The URL rewriter API is used by the crawler to filter and rewrite extracted URL links before they are inserted into the URL queue.

Ultra Search also includes highly functional query applications to query and display search results. The query applications are J2EE-compliant Web applications.

Ultra Search Query API

Ultra Search provides a Java API for querying indexed data. The API methods retrieve and display query results. Because it is written in Java, it is compatible with a large spectrum of Web application servers that support any Java-based technology, such as JSP version 1.1 and higher. The API uses JDBC connection pooling for scalability.

The Java API does not impose any HTML rendering elements. The application can completely customize the HTML interface. For example:

- Basic search form
- Advanced search form
- Query result display
- Help page
- Feedback page
- Register URL

You embed Ultra Search query functionality in your Web application with the supplied Ultra Search Java query API. The API supports two methods:

- Methods that retrieve query result data only.
- Methods that retrieve HTML code containing query result data.

The data-only methods do not return any HTML and can be used when you require full control over the HTML code to be rendered. The methods that retrieve HTML code support features such as allowing you to embed query input boxes and result lists in your Web application.

Some features of the Ultra Search Java query API include the following:

- Lets you retrieve query results
- Lets you set query properties, such as the total number of hits to return, and so on
- Lets you set the query session language
- Lets you access Ultra Search tables to retrieve Ultra Search dictionary data, such as all defined data groups and attributes
- Lets you customize and generate your query interface and search result screen with procedures that return blocks of HTML code that you can embed into your Web application
- Lets you allow the search end user to submit URLs to the seed URL list

The Ultra Search Java query API is encapsulated in the `oracle.ultrasearch.query` package.

See Also: "Tuning Query Performance" on page 4-3

Customizing the Query Syntax Expansion

Ultra Search uses the Oracle Text engine to index and search documents. When an end user specifies a certain query string, Ultra Search takes that string and transforms it into an Oracle Text query expression. This process is called query syntax expansion.

You can customize Ultra Search to use your own implementation of the query syntax expansion. In the Ultra Search 9.0.1 release, the default query syntax expansion implementation was contained in the `WK_QUERYEXP PL/SQL` package.

See Also: Appendix C, "Customizing the Query Syntax Expansion 9.0.1"

The default query expansion lets you specify a query syntax similar to most internet search engines. The syntax boosts scores for documents that match the user's query in the document 'title' string attribute. The syntax for Contains is the same when used on the document content and on string attributes.

The default query syntax expansion is implemented in the `oracle.ultrasearch.query.Contains` class. To customize query expansion, use the `oracle.ultrasearch.query.CtxContains` class.

This section describes the default query expansion rules, and how to customize the query syntax expansion to suit your organization's preferences.

Default Query Syntax Expansion Implementation

The default query syntax expansion implementation directly affects the following:

- The way the end user enters a query string (known as the *end user query syntax*)
- The way the documents matching the query are scored (known as *scoring*)
- The way the end user's query string is transformed into an Oracle Text query string (known as the *expansion rules*)

The default query syntax expansion is implemented in the `oracle.ultrasearch.query.Contains` class. The sample query applications makes use of this syntax expansion for content search as well as string attribute search.

End User Query Syntax

The end user query syntax defined by the default query syntax expansion implementation is similar to the standard text query syntax employed by most search engines on the Web.

- **Token:** A token is a string enclosed in double-quotes (""). It can be a single word or a phrase.
- **Operators:** The default implementation defines three operators. They are the [+], [-] and [*] operators. These operators are defined by the default implementation. Change these operators to whatever you prefer in your own custom implementation.

The plus operator [+] specifies that the token immediately following it must appear in all documents included in the search result.

The minus operator [-] specifies that the token immediately following it cannot appear in any document included in the search result.

The asterisk [*] specifies a wildcard search. It matches zero or more characters. A token starting with the asterisk is ignored. The asterisk can only be specified at the end (right side) or middle of a token. For example, "hel*o" and "hell*" use the asterisk correctly, but "*ello" is unacceptable.

The following table summarizes the rules for the Ultra Search end user query syntax:

Note: All end-user query strings are encased in square braces. For example, the end user query string Oracle Applications is notated as [Oracle Applications].

Rule	Description
Single word search	Entering one word finds documents that contain that word. For example, searching for [Oracle] finds all documents that contain the word "Oracle" anywhere in that document. Note: Searching for [Oracle] is not equivalent to [Oracle*].
Multiple word search	Entering more than one word finds documents that each contain any of those words in any order. For example, searching for [Oracle Applications] finds documents that contain "Oracle" or "Applications" or "Oracle Applications."
Compulsory inclusion [+]	Attaching a [+] in front of a word requires that the word be found in all matching documents. For example, searching for [Oracle + Applications] only finds documents that contain the word "Applications." Note: In a multiple word search, you can attach a [+] in front of every token including the very first token.
Compulsory exclusion [-]	Attaching a [-] in front of a word requires that the word must not be found in all matching documents. For example, searching for [Oracle - Applications] only finds documents that do not contain the word "Applications". Note: In a multiple word search, you can attach a [-] in front of every token except the very first token.
Phrase matching ["..."]	Putting quotes around a set of words only finds documents that contain that precise phrase. For example, searching for ["Oracle Applications"] finds only documents that contain the string "Oracle Applications."

Rule	Description
Wildcard matching [*]	Attaching a [*] to the right-hand side of a word returns left side partial matches. For example, searching for the string [Ora*] finds documents that contain all words beginning with "Ora," such as "Oracle" and "Orator." You can also insert an asterisk in the middle of a word. For example, searching for the string [A*e] retrieves documents that contain words such as "Apple", "Ate", "Ape", and so on. Wildcard matching requires more computational processing power and is generally slower than other types of queries.

Scoring Classes

There are three ways documents are matched against an end user query string. These three ways are known as scoring "classes." Documents are scored and ranked higher if they satisfy the requirements for a higher class. Within each class, documents are also ranked differently depending on how well they match the conditions of that scoring class.

Class 1 is the most heavily weighted class. The score is derived from the number of occurrences of a precise phrase in a document. A document that has more instances of the precise phrase have a higher score than another document that has fewer occurrences of the precise phrase.

Class 2 is the next more heavily weighted class. In this class, the closer the tokens appear in a document, the higher the score becomes. For example, an end user query string [Oracle Applications Financials] can result in three documents found. None of the three documents contain the precise phrase "Oracle Applications Financials." However, document X contains the all three tokens "Oracle", "Applications", and "Financials" in the same sentence separated by other words. Document Y contains the individual tokens in the same paragraph but in different sentences. Document Z contains the same three tokens, but each token resides in different paragraphs. In this scenario, document X has the highest score, because the tokens are closest together. Likewise, Y has a higher score than Z.

Class 3 is the least weighted class. A document that has more tokens gets a higher score. For example, an end user query string [Oracle Applications Financials] can result in three documents found. Document X might contain all three tokens. Document Y might contain the tokens "Oracle" and "Applications" only. Document Z might contain only the token "Oracle." In this scenario, document X has a higher score than Y. Likewise, Y has a higher score than Z.

Expansion Rules

As mentioned previously, the end user query is expanded to an Oracle Text query. The expanded query string rules are captured in BNF (Backus Naur Form) notation. Again, these rules are the rules that Ultra Search uses as a default query syntax expansion implementation.

The rules that define an expanded query:

```
<expanded query> ::= (<expression> within <title section>)*2, <expression>
```

```
<expression> ::= <generic query expression> | <simple query expression>
```

```
<generic query expression> ::= (([ <plus expression>*100 & ] (<main expression>))
[ <minus expression> ]
```

```
<simple query expression> ::= (<phrase expression>)*2, (<main expression>)
```

```
<main expression> ::= (<near expression>)*2, (<accum expression>)
```

The following list contains some terms and their meanings, which explain some of the terms used in the preceding rules:

A <plus expression> is an AND expression of all plus tokens.

A <minus expression> is a NOT expression of all minus tokens.

A <phrase expression> is a PHRASE formed by all tokens in the <main expression>

A <near expression> is a NEAR expression of all tokens but minus tokens.

An <accum expression> is an ACCUMULATE expression of all tokens but minus tokens.

A <simple query expression> is used only when the end user query has multiple tokens and does not have any operator or a double quote. Otherwise, a <generic query expression> is used.

If there is no token that is neither plus token nor minus token, then the <plus expression> and the <accum expression> are eliminated.

Examples of Applying the Rules

The following table illustrates how the default query syntax expansion implementation converts end user query strings into Oracle Text compatible query strings.

End User Query String	Expanded Query String Understandable by Oracle Text
[Oracle]	((({Oracle}) within TITLE__31)*2,({Oracle}))
[Oracle + Applications]	(((((Applications))*10)*10&({Oracle};{Applications})*2,({Oracle},{Applications})) within TITLE__31)*2,(((Applications))*10)*10&({Oracle};{Applications})*2,({Oracle},{Applications}))
[Oracle - Applications]	((({Oracle})~{Applications}) within TITLE__31)*2,(({Oracle})~{Applications}))
["Oracle Applications"]	((({Oracle Applications}) within TITLE__31)*2,({Oracle Applications}))
[Ora*]	((({Ora%}) within TITLE__31)*2,({Ora%}))
[Oracle Applications]	(((((Oracle Applications))*2,(({Oracle};{Applications})*2,({Oracle},{Applications})) within TITLE__31)*2,(((Oracle Applications))*2,(({Oracle};{Applications})*2,({Oracle},{Applications})))

Customizing the Rules

Customize this expansion to suit your organization's purposes by defining and implementing your own query syntax expansion. You should have detailed understanding of Oracle Text queries using the `ctxsys.contains()` operator. Oracle Text offers a rich set of linguistic features, such as thesaurus, theme, stemming, and soundex as a part of its query language.

See Also:

- *Oracle Text Application Developer's Guide*
- *Oracle Text Reference*

To customize Ultra Search to use your own implementation of the query syntax expansion, use the `oracle.ultrasearch.CtxContains` class in your query application instead of the `oracle.ultrasearch.query.Contains` class. `CtxContains` lets you use any Oracle Text query as a part of an Ultra Search query. Use the following steps:

1. Construct a Oracle Text query based on the user's input. For example, if the user's input is "cat", using the stemming feature, you can construct a Text query "\$cat", which will find documents with "cat" or "cats". You can use any tool to construct the Text query, as long as it is a string object. Depending on the complexity of user's query syntax, you might want to leverage some existing lexers in Java.
2. Construct a CtxContains using the Text query. For example:

```
String textQuery = "$cat";
oracle.ultrasearch.Query query = new oracle.ultrasearch.CtxContains
(textQuery);
```

The above code constructs a query for documents with "cat" or "cats". You can also limit that query to document titles (not content) as follows:

```
String textQuery = "cat";
StringAttribute titleAttribute =
instanceMetaData.getStringAttribute("TITLE");
oracle.ultrasearch.Query query = new oracle.ultrasearch.CtxContains
(textQuery, titleAttribute);
```

3. You can optionally combine the CtxContains with any other Ultra Search query by joining them with the And/Or query operators.
4. Run the query by invoking the getResult() method with the constructed query object.

See Also: Oracle Ultra Search Java API Reference for detailed information on the oracle.ultrasearch.query.CtxContains API

Ultra Search Query Tag Library

On top of the Java query API, Ultra Search provides a JSP tag library as an alternative for developing search applications. Based on the Sun Microsystems JavaServer Pages specification version 1.1, the Ultra Search tag library better separates the dynamic/Java development effort from the static/HTML development effort, and enables Web developers who are unfamiliar with Java to incorporate search functionality into their applications.

The Ultra Search tag library provides a subset of the features in the Java Query API. Advanced features, such as custom query expansion and URL submission, are not available as tags. The main features of the tag library are the following: ability to

retrieve search attributes, groups, languages, and LOVs for rendering the advance query form; and ability to iterate through the resulting hit set, and retrieve document attributes and properties for rendering the result page.

The tag library is summarized in following table:

Tag	Description	Attributes
instance	This tag establishes a connection to an Ultra Search instance.	instanceId username password URL dataSourceName tablePagePath emailPagePath filePagePath
showAttributes	For an advanced query, use this tag to show the list of attributes available.	instance locale
showGroups	For an advanced query, use this tag to show the list of groups.	instance locale
showLanguages	For an advanced query, use this tag to show the list of languages defined in the instance.	instance
showLOV	Show all values defined for a search attribute.	instance locale attributeName attributeType
getResult	Perform the search.	resultId instance query queryLocale documentLanguage from to boostTerm withCount

Tag	Description	Attributes
fetchAttribute	This is a nested tag within getResult to specify which attributes of each document should be fetched along with the query results. There can be any number of nested fetchAttribute tags.	attributeName attributeType
showHitCount	If withCount="true" in the getResult tag, then the result includes a total number of hits, and you can use showHitCount to display this number.	result
showResults	Renders the results of the search.	result instance
showAttributeValue	Renders a document attribute.	attributeName attributeType

Details of these tags are described in the following subsections. Note the following requirements for using Ultra Search tags:

- Install the file `ultrasearch_query.jar` and include it in classpath or the `WEB-INF/lib` directory of the Web application. This file is provided with the Ultra Search installation under the `ultrasearch/lib` directory.
- Make sure that the tag library description file, `ultrasearch-taglib.tld`, is deployed with the application and is in the location specified in the `taglib` directives of your JSP pages, such as in the following example: `<%@ taglib uri="/WEB-INF/ultrasearch-taglib.tld" prefix="US" %>`

The Ultra Search tag library definition (TLD) file can be found in `$ORACLE_HOME/ultrasearch/sample/query/WEB-INF/ultrasearch-taglib.tld` after `sample.ear` has been deployed. It is also packaged with `ultrasearch_query.jar` under the name `META-INF/taglib.tld`.

Query Tag Descriptions

The following section describes each Ultra Search tag, its attributes, and action. Examples are shown without any static HTML, which can be inserted to format the output.

<instance> Tag: Connecting to the Ultra Search Instance

This tag establishes a connection to an Ultra Search instance. Some basic parameters must be established for this tag to work, such as JDBC connection string, schema user name/password, Ultra Search instance name, and so on.

Attribute Name	Description
instanceId="name"	Names the instance defined by this tag. This name is then used by other Ultra Search tags to specify the instance being searched.
username	Creates a database connection.
password	Creates a database connection.
url	Gets the URL used to create a JDBC connection. This attribute is optional if dataSourceName is specified.
dataSourceName	The JNDI name that identifies a JDBC data source. Users should set either the URL or data source name properties. This is optional if URL is specified.
instanceName	The name of the Ultra Search instance that is owned by the schema user. If the schema user owns only one Ultra Search instance, then this is optional.
tablePagePath	The URL path of the Web application that renders the contents of a database table.
emailPagePath	The URL path of the Web application that renders the contents of an email.
filePagePath	The URL path of the Web application that renders the contents of a file.

This tag defines a scripting variable of the name set by the instanceId property. All the other tag properties correspond to a property in the `oracle.ultrasearch.query.QueryInstance` class. Either the URL or the `dataSourceName` attribute should be set. They are exclusive of each other.

The following example uses the URL property to connect to the database.

```
<US:instance
  instanceId="mybookstore"
  url="oracle:jdbc:thin:@dbhost:1521:inst1"
  username="scott"
  password="tiger"
  tablePage="../display.jsp"
  emailPage="../mail.jsp"
  filePage="../display.jsp"
/>
```


<iterAttributes> Tag: Show All Search Attributes

When a user wants to perform an advanced query, the application needs to show the list of attributes that are available, the list of groups, and the list of languages defined in the instance. This can be done using some iteration tags that define script variables for page rendering.

Each attribute in Ultra Search has a name, a type, and a display name that is translated depending on the locale that is set for the QueryInstance tag. The attribute type should be used to determine which operators can be used on this attribute and how to parse the user's input.

Attribute Name	Description
instance="name"	This is a mandatory attribute to refer to the object defined by the instance tag.
locale="locale"	This determines the display name fetched using this tag.

This tag is an iteration tag. It loops through all the search attributes in the instance referred to by the instance tag attribute. In each loop, it defines a scripting variable named "attribute", which is an `oracle.ultrasearch.query.Attribute` object. It also defines a string variable named "displayname", which is the localized name of the attribute.

The following example shows all the attributes in "mybookstore" instance, using their English display names.

```
<US:iterAttributes instance="mybookstore" locale="<%=Locale.ENGLISH%" >
<%= attribute %>
<%= displayname %>
</US:iterAttributes>
```

<iterGroups> Tag: Show All Search Groups

Similar to the showAttributes tag, the showGroups tag iterates through all the groups defined in an instance.

Attribute Name	Description
instance="name"	This is a mandatory attribute to refer to the object defined by the instance tag.
locale="locale"	This determines the display name fetched using this tag.

This tag loops through all the search groups in the instance referred to by the instance tag attribute. In each loop, it defines a scripting variable named "group", which is an `oracle.ultrasearch.query.Group` object. It also defines a string variable named "displayname", which is the localized name of the group.

The following example shows all the groups in "mybookstore" instance, using their English display names.

```
<US:iterGroups instance="mybookstore" locale="<%=Locale.ENGLISH%>" >
<%= group %>
<%= displayname %>
</US:iterGroups >
```

<iterLanguages> Tag: Show All Search Languages

Similar to the `showAttributes` tag, the `showLanguages` tag iterates through all the languages defined in an instance. Because each language is defined by a `java.util.Locale` object, their display names are not handled by Ultra Search. Therefore, this tag does not define the `displayname` scripting variable.

Attribute Name	Description
<code>instance="name"</code>	This is a mandatory attribute to refer to the object defined by the instance tag.

This tag is an iteration tag. It loops through all the search languages in the instance referred to by the instance tag attribute. In each loop, it defines a scripting variable named "language", which is a `java.util.Locale` object. The display name for the language is provided by Java as a property of the object itself (through the `getDisplayname()` method).

The following example shows all the languages in "mybookstore" instance, using their English display names.

```
<US:iterLanguages instance="mybookstore">
<%= language %>
<%= language.getDisplayName (Locale.ENGLISH) %>
</US:iterLanguages >
```

<iterLOV> Tag: Show All Values Defined for a Search Attribute

Attribute Name	Description
instance="name"	This is a mandatory attribute to refer to the object defined by the instance tag.
locale="locale"	This determines the display name fetched using this tag.
attributeName="attname"	The name of the attribute whose LOV is being fetched in this LOV.
attributeType="string number date"	The type of the attribute whose LOV is being fetched in this LOV. This is needed because attribute name does not uniquely identify an attribute in the instance.

This tag is an iteration tag. It loops through all the values in a search attribute's LOV. In each loop, it defines a scripting variable named "value", which is either a `java.lang.String`, `java.util.Date`, or `java.math.BigDecimal` object, depending on the attribute type. It also defines a string variable named "displayname", which is the localized display name of the value.

The following example shows all the values for a string attribute named "Dept" in "mybookstore" instance, using their English display names.

```
<US:iterLOV instance="mybookstore" attribute_name="Dept" attribute_type="String"
>
<%= value %>
<%= displayname %>
</US:iterLOV >
```

Formulating the Query

Ultra Search supports a set of classes for building queries. Currently these classes do not have any tag equivalents.

<getResult> Tag: Perform Search

This tag performs the search and returns the result by defining a scripting variable of the type `oracle.ultrasearch.query.Result`.

Attribute Name	Description
resultId="name"	This names the result generated by this tag. This name is then used by other tags to render the result on the page.

Attribute Name	Description
instance="name"	This is a mandatory attribute to refer to the object defined by the instance tag.
query="<%= expression %>"	This specifies a query object to search with.
queryLocale="locale"	This specifies the locale of the query object.
documentLanguage="locale"	This specifies the language of the documents for which to search. This is optional. If it is not specified, then all languages are included in the search.
from="number"	This specifies the index of the first hit.
to="number"	This specifies the index of the last hit.
boostTerm="string"	This specifies the search term that is used for relevance boosting. This is optional.
withCount="true false"	This specifies whether the result has an estimate of the total hit count. This is optional. If unspecified, the behavior is same as withCount=false.

The `<getResult>` tag corresponds to the `getResult()` method on the `oracle.ultrasearch.query.Instance` class. The attributes of tag map to the parameters of the method, with the exception that `getResult()` method can specify the attributes to fetch. The `<getResult>` tag require the use of the nested `<fetchAttribute>` tag to accomplish metadata selection.

The following example shows a search for the first 20 documents of a query in English that appears in French documents.

```
<US:getResult
  resultId="searchresult"
  instance="mybookstore"
  query=""
  queryLocale=""
  documentLanguage=""
  from="1" to="20">
</US:getResult>
```

<fetchAttribute> Tag: Metadata Selection

This tag is used as nested tag inside `<getResult>`. It specifies which attributes of each document should be fetched along with the query result. Each `<getResult>` can have any number of nested `<fetchAttribute>` tags.

Attribute Name	Description
attributeName="atname"	The name of the attribute whose LOV is being fetched in this LOV.
attributeType="string number date"	The type of the attribute whose LOV is being fetched in this LOV. This is needed because attribute name does not uniquely identify an attribute in the instance.

Each occurrence of the <fetchAttribute> adds to the list of attributes passed to the getResult() invoked by the <getResult> tag.

The following example shows the same search in <getResult> tag, but fetching title and publication-date attributes of each book.

```
<US:getResult
  resultId="searchresult"
  instance="mybookstore"
  query=""
  queryLocale=""
  documentLanguage=""
  from="1" to="20">
<US:fetchAttribute
  attributeName="title"
  attributeType="string" />
<US:fetchAttribute
  attributeName="publication-date"
  attributeType="date" />
</US:getResult>
```

<showHitCount> Tag: Show Estimated Hit Count

After the search is performed, the result must be rendered. If withCount=true is in the <US:getResult> tag, then the result contains a count of total hits, and <showHitCount> tag can be used to display it.

Attribute Name	Description
result="name"	This refers to the resultId specified in the <US:getResult> tag.

This tag outputs the hit count to the page.

The following shows the hit count of the a search result.

```
<US:showHitCount result="searchresult" />
```

<iterResult> Tag: Render the Results

This tag is an iteration tag. It loops through all the documents in a search result.

Attribute Name	Description
result="name"	This refers to the resultId specified in the <US:getResult> tag.
instance="name"	This refers to the instanceId specified in the <US:instance> tag.

The tag loops through all the documents in a search result and defines a scripting variable "doc" that is a `oracle.ultrasearch.query.Document` object. In addition, it can have nested tags of <showAttributeValue>, which helps to render the document's attributes. It is an error if the result specified is not one obtained from search on the instance specified. In other words, the result must come from the instance.

The following example shows the URL of all documents in a search result.

```
<US:iterResult
result="searchresult"
instance="mybookstore">
</US:iterResult>
```

<showAttributeValue> Tag: Render a Document Attribute

This tag shows an attribute of a document within the <US:iterResult> tag.

Attribute Name	Description
attributeName="attname"	The name of the document attribute.
attributeType="string number date"	The type of the document attribute. This is needed because attribute name does not uniquely identify an attribute in the instance.
default="default string"	A value to output when the document has no value for this attribute. This is useful when a document has no title. The string "No Title" can be displayed as the default value.

This tag looks up the document attribute value and renders it on the page. If the attribute was not fetched as part of the search result, then nothing is output to the page.

The following example shows the title and publication dates of all documents in a search result.

```
<US:iterResult
result="searchresult"
instance="mybookstore">
<US:showAttributeValue attributeName="title" attributeType="string" default="No
Title" />
<US:showAttributeValue attributeName="publication-date" attributeType="date" />
</US:iterResult>
```

Ultra Search Crawler Agent API

You can implement a crawler agent to crawl and index a proprietary document repository, such as Lotus Notes or Documentum. In Ultra Search, the proprietary repository is called a user-defined data source. The module that enables the crawler to access the data source is called a crawler agent.

The agent collects document URLs and associated metadata from the user-defined data source and returns the information to the Ultra Search crawler, which enqueues it for later crawling. The crawler agent must be implemented in Java using the Ultra Search crawler agent API.

Ultra Search provides a sample implementation of user-defined crawler agents using the Ultra Search agent API. Upon invocation, this sample agent connects to a specified Oracle database and retrieves the contents of a table for the crawler to collect and index.

The sample agents are fully functional and can be customized to adapt to other database-based data sources. These agents performs the following tasks:

- Read data source parameters
- Connect to the database that contains the data source
- Initialize fetching document URL and attributes from the data source
- Fetch document URL and attributes from the data source
- Disconnect from the data source

Crawler Agent Overview

A crawler agent does the following:

- Authenticates the crawler for accessing the data source

- Provides access to the data source document through a HTTP URL (display URL)
- Provides the metadata of the document in the form of document attributes
- Maps each document attribute to a common attribute name used by end users
- Provides a "flattened" view of the data source, such that documents are retrieved one by one in a streaming fashion
- Instructs the crawler to parse the URL document for standard metadata, like author and title, if necessary
- Optionally provides the list of URLs that have changed since a given time stamp
- Optionally provides an access URL in addition to the display URL for the processing of the document

From the crawler's perspective, the agent retrieves the list of URLs from the target data source and saves it in the crawler queue before processing it.

Note: If the crawler is interrupted for any reason, then the agent invocation process is repeated with the original last crawl time stamp. If the crawler finished enqueueing URLs fetched from the agent and is half way done crawling, then the crawler only starts the agent, but does not try to fetch URLs from the agent. Instead, it finishes crawling the URLs already enqueued.

There are two kinds of crawler agents:

- Standard Agent
- Smart Agent

Standard Agent

The standard agent returns the list of URLs currently existing in the data source. It does not know whether any of the URLs had been crawled before, and it relies on the crawler to find any updates to the target data source. The standard agent's interaction with the crawler is the following:

- Crawler marks all existing URLs of this data source for garbage collection, assuming they no longer exist in the target data source.

- Crawler calls the agent to get an updated list of URLs. It marks for crawling every URL that already exists. If it is new, it inserts it into the URL table and queue.
- Crawler deletes the URLs that are still marked for garbage collection.
- Crawler goes through every URL marked for crawling and checks for updates.

Smart Agent

The smart agent uses a modified-since time stamp (provided by the crawler) to return the list of URLs that have been updated, inserted, and deleted. The crawler only crawls URLs returned by the agent and does not recrawl existing ones. For URLs that were deleted, the crawler removes them from the URL table. If the smart agent can only return updated or inserted URLs but not deleted URLs, then deleted URLs are not detected by the crawler. In this case, you must change the schedule crawler recrawl policy to periodically run the schedule in force recrawl mode. Force recrawl mode signals to the agent to return every URL in the data source.

The agent API `isDeltaCrawlingCapable()` tells the crawler whether the agent it invokes is a standard agent or a smart agent. The agent API `startCrawling(boolean forceRecrawl, Date lastCrawlTime)` lets the crawler tell the agent the last crawl time and whether the crawler is running in force recrawl mode.

Document Attributes and Properties

Document attributes, or metadata, describe document properties. Some attributes can be irrelevant to your application. The crawler agent creator must decide which document attributes should be extracted and saved. The agent also can be created such that the list of collected attributes are configurable. Ultra search automatically registers attributes returned by the agent. The agent can decide which attributes to return for a document.

Crawler Agent Functionality

This section describes aspects of the crawler agent.

Data Source Type Registration

A data source type is an abstraction of a data source. You can define new data source types with the following attributes:

- Name of data source type: For example, Lotus Notes. The name cannot be more than 100 bytes.
- ID of data source type: This is automatically assigned.
- Description of the data source type: This limit is 4000 bytes.
- Agent Java class name: For example, WebDbAgent. The location of this class is predefined by Ultra Search in `$ORACLE_HOME/ultrasearch/lib/agent/` and cannot be changed.
- Agent Java jar file name: The agent class can be stored in a Java jar file. This jar file must be in `$ORACLE_HOME/ultrasearch/lib/agent/`, where `$ORACLE_HOME` is the Oracle home directory where the Ultra Search backend (server component), *not* the middle tier, is installed.
- Parameters: Parameters are the properties of a data source; for example, seed URL, inclusion pattern, and robots exclusion for a Web data source. Define a parameter by specifying a parameter name (100 bytes maximum) and a description (4000 bytes maximum). By default, a parameter is not encrypted.
- Encryption: Should the value of this parameter be encrypted when stored.

Ultra Search does not enforce the occurrence of parameters. You cannot specify a particular parameter to have 0 or more, at least 1, or only 1 occurrence.

Data Source Registration

After a data source type is defined, any instance of that data source type can be defined:

- Data source name
- Description of the data source; limit to 4000 bytes
- Data source type ID
- Default language; default is 'en' (English)
- Parameter values; for example, seed - `http://www.oracle.com depth - 8`

Data Source Attribute Registration

You can add new attributes to Ultra Search by providing the attribute name and the attribute data type. The data type can be string, number, or date. Attributes with the same name but different data type can be added. Attributes returned by an agent are automatically registered if they have not been defined.

User-Implemented Crawler Agent

The crawler agent has the following requirements:

- The agent must be implemented in Java.
- The agent must support the Java agent APIs defined by Ultra Search.
- The agent must return the URL attributes and properties.
- The agent optionally can authenticate the crawler's access to the data source.
- The agent must "flatten" the data source such that each document is retrieved one by one in a streaming fashion. This is to encapsulate the crawling logic of a specific data source into the agent.
- The agent must decide which document attributes Ultra Search should keep. Any attribute not defined in Ultra Search is registered automatically.
- The agent can map attributes to data source properties. For example, if an attribute "ID" is the unique ID of a document, then the agent should return (document_key, 4) where "ID" has been mapped to the property "document_key" and its value is 4 for this particular document.
- If the attribute LOV is available, then the agent returns them upon request.

Interaction Between the Crawler and the Crawler Agent

The crawler crawls data sources defined by the user through the invocation of the user-supplied crawler agent. The crawler can do the following:

- Invoke the crawler agent of the defined data source
- Supply data source parameter information to the agent
- Authenticate itself with the agent if needed
- Retrieve a list of URLs and associate attributes/properties that must be crawled
- Use the URL provided by the agent to retrieve the document
- Detect insert, update, and delete to the data source
- Retrieve attribute LOV data if available

Crawler Agent APIs and Classes

The crawler agent API is a collection of methods used to implement a crawler agent. A sample implementation of a crawler agent `SampleAgent.java` is provided under `$ORACLE_HOME/ultrasearch/extension/`.

UrlData: The crawler agent uses this interface to populate document properties and attribute values. Ultra Search provides a basic implementation of this interface that the agent can use directly or extend if necessary. The class is `DocAttributes` with a constructor that has no argument. The agent might decide to create a pool of `UrlData` objects and cycle through them during crawling. In the most simple implementation, the agent creates one `DocAttributes` object, repeatedly resets and populates the data, and returns this object.

LovInfo: The crawler agent uses this interface to submit attribute LOV definitions.

DataSourceParams: The crawler agent uses this interface to read and write data source parameters.

AgentException: The crawler agent uses this exception class when an error occurs.

CrawlerAgent: This interface lets the crawler communicate with the user-defined data source. The crawler agent must implement this interface.

Sample Agent Files

The sample agent files are located in the `$ORACLE_HOME/ultrasearch/extension` directory. You can view the sample agent source code using your preferred text editor.

There is a `SampleAgent_readme.htm` file and a `SampleAgent.java` file. These are for the sample crawler agent implementation using agent APIs.

Setting up the Sample Crawler Agent

This section describes how to set up the sample crawler agent.

Compiling and Building the Agent Jar File

The Java source code for the sample agent first must be compiled into class files and put into a jar file in the `$ORACLE_HOME/ultrasearch/lib/agent/` directory, where `$ORACLE_HOME` is the Oracle home directory where the Ultra Search backend (server component), *not* the middle tier, is installed.

The classes needed for compilation are the JDK class (`classes.zip`), Oracle JDBC Thin Driver (`classes12.zip`), and `ultrasearch.jar`. For example:

```
javac -J-ms16m -J-mx96m -O -classpath /jdk1.2.2
05/lib/classes.zip:/lib/classes12.zip:
$ORACLE_HOME/ultrasearch/lib/ultrasearch.jar SampleAgent.java
```

To build the `SampleAgent.jar` file, enter the following:

```
/jdk1.2.2_05/bin/jar cv0f /oracle/ultrasearch/lib/agent/SampleAgent.jar  
SampleAgent.class 'SampleAgent$DocNode.class'
```

Creating a Data Source Type

A data source type that uses the sample agent must be created first.

- Name: URL table type
- Description: Table with rows of URLs
- Agent Name: SampleAgent
- Agent Jar File: sampleagent

Defining Data Source Parameters

Define parameters for a data source type:

- Database Connect String (DB connection)
- User Name (schema owner of the URL table)
- Password (schema owner password, encrypted)
- Table Name (URL table name)
- URL Column (Column holding doc URLs)
- Ignore Flag Column (1 for ignoring, 0 otherwise)
- Language Column (Document Language)
- Attribute List (List of column for attributes)
- It is in the following format: [column name/attribute name] <data type> [column name/attribute name] <data type> ... where <data type> 0 is number, 1 is string, and 2 is date. For example, if the document has 4 attributes: Company Name, Category, Revenue, S&P Rating, then it is specified as: [Company Name/Company/1][Category/Classification/1][Revenue/Revenue/0][Rating/Analyst Rating/1]
- Log File Name (log file)
- Log Directory (Location of log file)

Defining a Data Source of this Type

A data source is defined, which initializes the data source parameters. For example, the value specified accesses a table whose schema is the following:

```
TABLE NEWS (  
  ARTICLE_NO    NUMBER,  
  NEWS_URL      VARCHAR2(740),  
  TITLE         VARCHAR2(200),  
  AUTHOR        VARCHAR2(100),  
  PUB_DATE      DATE default SYSDATE,  
  PUBLISHER     VARCHAR2(100),  
  PRICE         NUMBER,  
  LANG          VARCHAR2(10),  
  IGNORE       NUMBER DEFAULT 0,  
  PRIMARY KEY (NEWS_URL)  
);
```

- Database Connect String: dlsun1710:5521:search
- User Name: SCOTT
- Password: TIGER
- Table Name: NEWS
- URL Column: NEWS_URL
- Ignore Flag Column: IGNORE
- Language Column: LANG
- Attribute List: [ARTICLE_NO/Article Number/0][TITLE/Article Title/1][AUTHOR/Author/1][PUB_DATE/Report Date/2][PUBLISHER/Newspaper/1][PRICE/Download Cost/0]
- Log File Name: testagent.log
- Log Directory: /tmp/ultrasearch/

Ultra Search Java Email API

Ultra Search provides a Java API for accessing archived emails. The Ultra Search query application uses the API to display emails addressed to mailing lists that have been indexed by the Ultra Search system. The API can also be used to build your own custom query application.

The application user-interface logic is entirely controlled in the JSP. Therefore, you can customize the look-and-feel to your needs.

Email documents contain valuable information, but they are not structured to find specific relevant information easily. Ultra Search lets you retrieve and index emails on a server that supports the IMAP4 protocol.

An email source is a data source that derives its content from emails sent to a specific email address. When the Ultra Search crawler searches an email source, the crawler collects all emails that have the specific email address in any of the "To:" or "Cc:" email header fields.

Note: Ultra Search stores copies of all retrieved emails in the local file system of the Ultra Search server installation.

A possible application of an email source is where an email source represents all emails sent to a mailing list. In such a scenario, multiple email sources are defined where each email source represents an email list.

Ultra Search email crawling and rendering is built on top of the JavaMail API using Sun Microsystems' reference implementation of JavaMail. This enables Ultra Search to provide a Java API for accessing indexed emails. The API is known as the Ultra Search Java Email API. This API lets you retrieve information such as email header information, email body content, and attachments of an email.

Use this API to embed Ultra Search email browsing functionality into JavaServer Page (JSP) or servlet-based Web applications. Ultra Search ships a fully functional JSP Web application that directly uses this API to render indexed emails. Because the source code is viewable, you can use it as an example for building your own customized email browser.

JavaMail Implementation

Ultra Search requires a JavaMail 1.1 compliant implementation. The reference implementation by Sun Microsystems is JavaMail version 1.2. This reference implementation is shipped with Ultra Search.

Java Email API

The Ultra Search Java Email API is encapsulated in the `oracle.ultrasearch.query` package.

Sample Mailing List Browser Application Files

The sample mailing list browser applications files are located in the `$ORACLE_HOME/ultrasearch/sample/query` directory. You can directly view the sample mailing list browser application source code using your preferred text editor.

The following tables describe all sample mailing list browser application files, README file, and stylesheets:

File	Description
SampleAgent_readme.html	Readme
mail.css	Style sheet for sample email Web application

Sample JavaServer Page Mailing List Browser Applications Files:

File	Description
mail.jsp	Mailing list browser applications that selectively include HTML code returned by other JSP files, depending on what the end user wants to view
mailindex.jsp	JSP page that displays all email sources (mailing lists) of an Ultra Search instance
mailmsgs.jsp	JSP page that displays all emails for an email source (mailing list)
mailreader.jsp	JSP page that displays an email
mailutil.jsp	JSP page that defines various functions that are used by mailreader.jsp

Graphics Files for All Applications:

File	Description
images/ultra_mediumbanner.gif	Ultra Search banner
images/wsd.gif	Background image used in sample query application

Setting up the Sample Mailing List Browser Application

For detailed instructions on setting up the sample JSP mailing list browser application, see "Installing the Ultra Search Middle Tier on Web Server Hosts" on page 2-11.

Ultra Search URL Rewriter API

A URL rewriter is a user supplied Java module that implements the Ultra Search `UrlRewriter` Java interface. When activated, it is used by the crawler to filter and rewrite extracted URL links before they are inserted into the URL queue.

Web crawling generally consists of the following steps:

1. Get the next URL from the URL queue. (Web crawling stops when the queue is empty.)
2. Fetch the contents of the URL.
3. Extract URL links from the contents.
4. Insert the links into the URL queue.

The generated new URL link is subject to all existing host, path, and mimetype inclusion and exclusion rules.

There are two possible operations that can be done on the extracted URL link:

- Filtering: removes the unwanted URL link
- Rewriting: transforms the URL link

URL Link Filtering

Users control what type of URL links are allowed to be inserted into the queue with the following mechanisms supported by the Ultra Search crawler:

- `robots.txt` file on the target Web site; for example, disallow URLs from the `/cgi` directory
- Hosts inclusion and exclusion rules; for example, only allow URLs from `www.acme.com`
- File path inclusion and exclusion rules; for example, only allow URLs under the `/archive` directory
- Mimetype inclusion rules; for example, only allow HTML and PDF files

- Robots metatag `NOFOLLOW`; for example, do not extract any link from that page
- Black list URL; for example, URL explicitly singled out not to be crawled

Note: All URLs must pass domain rules before being checked for path rules. Path rules let you further restrict the crawling space. Path rules are host-specific, but you can specify more than one path rule for each host. For example, on the same host, you can include path files: `//host/doc` and exclude path files: `//host/doc/unwanted`.

With these mechanisms, only URL links that meet the filtering criteria are processed. However, there are other criteria that users might want to use to filter URL links. For example:

- Allow URLs with certain file name extensions
- Allow URLs only from a particular port number
- Disallow any PDF file if it is from a particular directory

The possible criteria could be very large, which is why it is delegated to a user-implemented module that can be used by the crawler when evaluating an extracted URL link.

URL Link Rewriting

For some applications, due to security reasons, the URL crawled is different from the one seen by the end user. For example, crawling is done on an internal Web site behind a firewall without security checking, but when queried by an end user, a corresponding mirror URL outside the firewall must be used.

A *display URL* is a URL string used for search hit display. This is the URL used when users click the search hit link. An *access URL* is a URL string used by the crawler for crawling and indexing. An access URL is optional. If it does not exist, then the crawler uses the display URL for crawling and indexing. If it does exist, then it is used by the crawler instead of the display URL for crawling.

For regular Web crawling, there are only display URLs available. But in some situations, the crawler needs an access URL for crawling the internal site while keeping a display URL for the external use. For every internal URL, there is an external mirrored one.

For example:

`http://www.acme-ga.us.com:9393/index.html`
`http://www.acme.com/index.html`

When the URL link `http://www.acme-ga.us.com:9393/index.html` is extracted and before it is inserted into the queue, the crawler generates a new display URL and a new access URL for it:

Access URL:

`http://www.acme-ga.us.com:9393/index.html`

Display URL:

`http://www.acme.com/index.html`

The extracted URL link is rewritten, and the crawler crawls the internal Web site without exposing it to the end user.

Another example is when the links that the crawler picks up are generated dynamically and can be different (depending on referencing page or other factor) even though they all point to the same page. For example:

`http://compete3.acme.com/rt/rt.wvw_media.show?p_type=text&p_id=4424&p_currcornerid=281&p_textid=4423&p_language=us`

`http://compete3.acme.com/rt/rt.wvw_media.show?p_type=text&p_id=4424&p_currcornerid=498&p_textid=4423&p_language=us`

Because the crawler detects different URLs with the same contents only when there is sufficient number of duplication, the URL queue could grow to a huge number of URLs, causing excessive URL link generation. In this situation, allow "normalization" of the extracted links so that URLs pointing to the same page have the same URL. The algorithm for rewriting these URLs is application dependent and cannot be handled by the crawler in a generic way.

When a URL link goes through a rewriter, there are the following possible outcomes:

- The link is inserted with no changes made to it.
- The link is discarded; it is not inserted.
- A new display URL is returned, replacing the URL link for insertion.
- A display URL and an access URL are returned. The display URL may or may not be identical to the URL link.

Creating and Using a URL Rewriter

Follow these steps to create and use a URL rewriter:

1. Create a new Java file implementing the `UrlRewriter` interface `open()`, `close()`, and `rewrite()` methods. A sample rewriter, `SampleRewriter.java`, is available for reference under `$ORACLE_HOME/ultrasearch/extension/`.

2. Compile the rewriter Java file into a class file. For example:

```
/jdk1.3.1/bin/javac -O -classpath $ORACLE_
HOME/ultrasearch/lib/ultrasearch.jar SampleRewriter.java
```

3. Package the rewriter class file into a jar file under the `$ORACLE_HOME/ultrasearch/lib/agent/` directory. For example:

```
/jdk1.3.1/bin/jar cv0f $ORACLE_HOME/ultrasearch/lib/agent/sample.jar
SampleRewriter.class
```

4. Specify the rewriter class name and jar file name (for example, `SampleRewriter` and `sample.jar`) in the administration tool in step 2 of "Creating Web Sources" on page 7-21 or in the crawler parameters page of an existing Web data source.
5. Enable the `UrlRewriter` option from Web Sources page in the administration tool.
6. Crawl the target Web data source by launching the corresponding schedule. The crawler log file confirms the use of the URL rewriter with the message *Loading URL rewriter "SampleRewriter"...*

Note: URL rewriting is available for Web data sources only.

See Also:

- *Oracle Ultra Search API Reference* for the API (`oracle.ultrasearch.crawler` package)
- The sample URL rewriter `SampleRewriter.java` under `$ORACLE_HOME/ultrasearch/extension/`
- "Web Sources" on page 7-20

Ultra Search Sample Query Applications

Ultra Search provides several sample query applications and a sample crawler agent. Use the sample query applications as examples for creating your own query application. The query applications are written as JavaServer Page (JSP) applications. Your query application uses the Ultra Search query API. You can also use the sample crawler agent to create your own crawler agent.

Note: Pointers to the sample query applications and the sample crawler agent Java source code, as well as their corresponding readmes, are in the Ultra Search welcome page:

```
http://hostname.domainname:port/ultrasearch/index.html
```

The sample query applications are shipped as a deployed J2EE Web application (`sample.ear`). This component depends on a J2EE container to host the Web pages, a JDBC driver, and Java Mail API for displaying email results. After the `sample.ear` file is deployed by the Oracle Containers for J2EE (OC4J), you see a set of JSP files that demonstrate the query API usage.

The sample query applications include a sample search portlet. The sample Ultra Search portlet demonstrates how to write a search portlet for use in OracleAS Portal.

When the user issues a query in any of the query applications, a hit list containing query results is returned. The user can select a document to view from the hit list. A hit list can include HTML documents, files, database table content, archived emails, or OracleAS items. The Ultra Search sample query applications also incorporate an email browser for reading and browsing emails.

The Ultra Search administration tool and the Ultra Search sample query applications are part of the Ultra Search middle tier. However, the Ultra Search administration tool is independent from the Ultra Search sample query applications. Therefore, they can be hosted on different computers to enhance security or scalability.

If you do not want to use the sample query applications, you can build your own query application by directly invoking the Ultra Search Java Query API. Because the API is coded in Java, you can invoke the API methods from any Java-based application, such as from a Java servlet or a JavaServer Page (as in the case of the provided sample query applications). For rendering emails that have been crawled and indexed, you can also directly invoke the Ultra Search Java email API methods.

Sample Query Applications

The sample query applications are located in the `$ORACLE_HOME/ultrasearch/sample` directory.

JavaServer Page Concepts

As mentioned earlier, you can use JSP code and the supplied Java APIs to create your Web application. Typically, your Web application runs in an application server, such as Oracle Application Server. The application server typically runs on a separate computer from the Oracle server for performance and scalability reasons. The Oracle server holds the Ultra Search indexes.

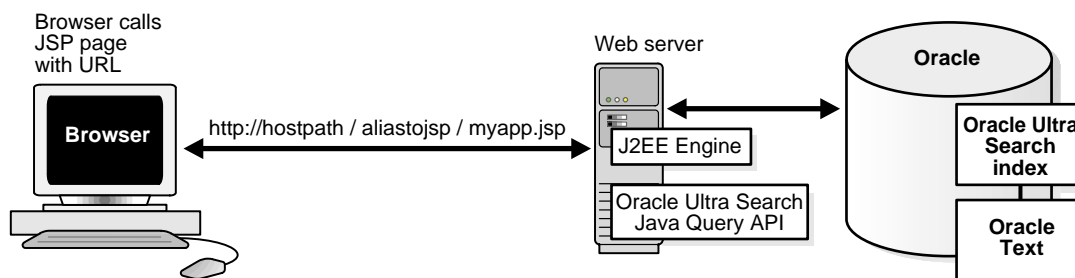
JSP applications are compiled into Java servlets at runtime. The compiled servlets run in one or more Java Virtual Machine processes. The JSP application communicates with the Oracle server through the Oracle JDBC driver.

As in any Java application, you must include the following files in your servlet engine classpath to use the Java query and email APIs:

- `$ORACLE_HOME/ultrasearch/lib/ultrasearch_query.jar`
- `$ORACLE_HOME/lib/mail.jar`
- `$ORACLE_HOME/lib/activation.jar`

Figure 8–1 shows how your Web query application calls the Ultra Search Java query API.

Figure 8–1 Calling JavaServer Pages



Loading Metadata into Ultra Search

Ultra Search provides a command-line tool to load metadata into an Ultra Search database. If you have a large amount of data, then this is probably faster than using the HTML-based administration tool.

The loader tool supports the following types of metadata:

- Search attribute list of values (LOVs) and display names
- Document relevancy boosting and document loading

The metadata loader is a Java application. To use the program, you must put the metadata in an XML file that conforms to the XML schema formats described in the following sections. You then can launch the Java program with the XML filename, the database related parameters, and the loader type parameter. The program parses the XML file and uploads the metadata. Status and error messages are displayed in the terminal console.

See Also: "Document Relevancy Boosting" on page 1-10

Launching the Loading Tool

The loader program binary file is located in the following directory:
`%ULTRASEARCH_HOME%/bin/MetaLoader.class.`

Your computer should have Java 1.2 compliant Java Runtime or higher. The following Java libraries should be included in the system Java CLASSPATH:

- Oracle JDBC Thin Driver version 1.2. The filename is `classes12.zip`.
- Oracle XML parser for Java version 2. The filename is `xmlparserv2.jar`.
- Oracle XML schema processor for Java. The filename is `xmlschema.jar`.
- Ultra Search Java library. The filename is `ultrasearch.jar`.

- Oracle JDBC globalization support version 1.2. The filename is `nls_charset12.zip`.

To launch the file, enter the following:

```
% java MetaLoader -db database_connection_string -u user_name -p password  
-i instance_name -type loader_type -f input_file
```

Where:

- `-db` is the database connection string
- `-u` is the database schema user name
- `-p` is the database schema password
- `-i` is the Ultra Search instance name
- `-type` is the loader metadata type:lov or doc
- `-f` is the input metadata XML filename

For example, suppose you use the tool to load attribute LOVs specified in the XML file `test.xml` with the following arguments:

- Database connection string: `dlsun576:5521:isearch`
- Schema user name: `wk_test`
- Schema password: `welcome`
- Ultra Search instance name: `wk_inst`

The following statement launches the loader program:

```
% java MetaLoader -db dlsun576:5521:isearch -u wk_test -p welcome -i wk_inst  
-type lov -f test.xml
```

Loading Documents and Relevance Scores

To use the loader tool to add documents and their relevancy boosting scores into Ultra Search, the parameter `-type` value should be `doc`.

The Input XML File

The document URL and relevance boosting scores are defined in an XML file. You can define one or more documents to be boosted. Each document can have one or more boosting score pairs. The definition of the XML file is stored in the XML schema.

See Also: "XML Schema for Document Relevance Boosting" on page A-5

Example of the Document Relevance Boosting XML File

```
<?xml version = "1.0" encoding = "UTF-8"?>
<doc_list>
  <doc url="http://www.oracle.com" data_source_name="Data Source A">
    <term score="100">database</term>
    <term score="90">internet</term>
    <term score="80">software</term>
  </doc>
  <doc url="http://www-st.us.oracle.com" data_source_name="Data Source B">
    <term score="100">Sever Technology</term>
    <term score="100">ST Web site</term>
    <term score="95">st</term>
  </doc>
</doc_list>
```

In the previous example, the document URL `http://www.oracle.com` is loaded to the data source `Data Source A`. This is defined in Ultra Search with relevance boosting term `database` and score 100, term `internet` and score 90, term `software` and score 80.

Note: The data source name is the original data source name, not the data source display name.

Loading Search Attribute LOVs and LOV Display Names

To use loader tool to add LOV entries and display names to Ultra Search, the parameter `-type` value should be `lov`.

The LOV XML File

The LOV entries and display names are defined in a XML file. You can define one or more search attribute LOVs in the XML file. Both default LOV and data source-specific LOVs are put in the XML file. The definition of the XML file is stored in the XML schema.

See Also: "XML Schema for LOVs and LOV Display Names" on page A-5

Example of the LOV XML File

```
<?xml version = "1.0" encoding = "UTF-8"?>
<lov_list>
  <lov search_attr_name="Department" search_attr_type="string">
    <default>
      <lov_values>
        <entry value="100"></entry>
        <entry value="200"></entry>
      </lov_values>
      <lov_display_names lang="en-US">
        <entry value="100" display_name="Human Resource"></entry>
        <entry value="200" display_name="Finance"></entry>
      </lov_display_names>
    </default>
    <data_source name = "data source a">
      <lov_values>
        <entry value="300"></entry>
        <entry value="400"></entry>
      </lov_values>
      <lov_display_names lang="en-US">
        <entry value="300" display_name="Sales"></entry>
        <entry value="400" display_name="Marketing"></entry>
      </lov_display_names>
    </data_source>
    <data_source name = "data source b">
      <lov_values>
        <entry value="500"></entry>
        <entry value="600"></entry>
      </lov_values>
      <lov_display_names lang="en-US">
        <entry value="500" display_name="Production"></entry>
        <entry value="600" display_name="Research"></entry>
      </lov_display_names>
    </data_source>
  </lov>
</lov_list>
```

In the previous example, several LOVs for the string type search attribute Department are loaded to Ultra Search. They are:

- Default LOV entries for search attribute Department
- Search attribute Department LOV for data source data source a
- Search attribute Department LOV for data source data source b

XML Schema for Document Relevance Boosting

The XML schema for document relevance boosting terms and scores is described as follows:

```
<?xml version = "1.0" encoding = "UTF-8"?>
<!--Generated by XML Authority. Conforms to w3c http://www.w3.org/2001/XMLSchema-->
<xsd:schema xmlns:xsd = "http://www.w3.org/2001/XMLSchema"
  elementFormDefault = "qualified">
  <xsd:element name = "doc_list">
    <xsd:complexType>
      <xsd:sequence>
        <xsd:element name = "doc" maxOccurs = "unbounded">
          <xsd:complexType>
            <xsd:sequence>
              <xsd:element name = "term" maxOccurs = "unbounded">
                <xsd:complexType>
                  <xsd:simpleContent>
                    <xsd:extension base = "xsd:string">
                      <xsd:attribute name = "score" use = "required" type = "xsd:integer"/>
                    </xsd:extension>
                  </xsd:simpleContent>
                </xsd:complexType>
              </xsd:element>
            </xsd:sequence>
          </xsd:complexType>
          <xsd:attribute name = "url" use = "required" type = "xsd:string"/>
          <xsd:attribute name = "data_source_name" use = "required" type = "xsd:string"/>
        </xsd:element>
      </xsd:sequence>
    </xsd:complexType>
  </xsd:element>
</xsd:schema>
```

XML Schema for LOVs and LOV Display Names

The XML schema for LOV entries and display names is described as follows:

```
<?xml version = "1.0" encoding = "UTF-8"?>
<!--Generated by XML Authority. Conforms to w3c http://www.w3.org/2001/XMLSchema-->
<xsd:schema xmlns:xsd = "http://www.w3.org/2001/XMLSchema"
  elementFormDefault = "qualified">
  <xsd:element name = "lov_list">
    <xsd:complexType>
      <xsd:sequence>
```

```

<xsd:element name = "lov" maxOccurs = "unbounded">
  <xsd:complexType>
    <xsd:sequence>
      <xsd:element name = "default" minOccurs = "0">
        <xsd:complexType>
          <xsd:sequence>
            <xsd:element name = "lov_values" minOccurs = "0">
              <xsd:complexType>
                <xsd:sequence>
                  <xsd:element name = "entry" maxOccurs = "unbounded">
                    <xsd:complexType>
                      <xsd:attribute name = "value" use = "required" type =
"xsd:string"/>
                      </xsd:complexType>
                    </xsd:element>
                  </xsd:sequence>
                </xsd:complexType>
              </xsd:element>
            <xsd:element name = "lov_display_names" minOccurs = "0" maxOccurs =
"unbounded">
              <xsd:complexType>
                <xsd:sequence>
                  <xsd:element name = "entry" maxOccurs = "unbounded">
                    <xsd:complexType>
                      <xsd:attribute name = "value" use = "required" type =
"xsd:string"/>
                      <xsd:attribute name = "display_name" use = "required" type =
"xsd:string"/>
                      </xsd:complexType>
                    </xsd:element>
                  </xsd:sequence>
                  <xsd:attribute name = "lang" use = "required">
                    <xsd:simpleType>
                      <xsd:restriction base = "xsd:string">
                        <xsd:length value = "5"/>
                        <xsd:pattern value = "[a-zA-Z]{2}\-[a-zA-Z]{2}"/>
                      </xsd:restriction>
                    </xsd:simpleType>
                  </xsd:attribute>
                </xsd:complexType>
              </xsd:element>
            </xsd:sequence>
          </xsd:complexType>
        </xsd:element>
      <xsd:element name = "data_source" minOccurs = "0" maxOccurs = "unbounded">

```

```

<xsd:complexType>
  <xsd:sequence>
    <xsd:element name = "lov_values" minOccurs = "0">
      <xsd:complexType>
        <xsd:sequence>
          <xsd:element name = "entry" maxOccurs = "unbounded">
            <xsd:complexType>
              <xsd:attribute name = "value" use = "required" type =
"xsd:string"/>
              </xsd:complexType>
            </xsd:element>
          </xsd:sequence>
        </xsd:complexType>
      </xsd:element>
    <xsd:element name = "lov_display_names" minOccurs = "0">
      <xsd:complexType>
        <xsd:sequence>
          <xsd:element name = "entry" maxOccurs = "unbounded">
            <xsd:complexType>
              <xsd:attribute name = "value" use = "required" type =
"xsd:string"/>
              <xsd:attribute name = "display_name" use = "required" type =
"xsd:string"/>
              </xsd:complexType>
            </xsd:element>
          </xsd:sequence>
          <xsd:attribute name = "lang" use = "required">
            <xsd:simpleType>
              <xsd:restriction base = "xsd:string">
                <xsd:length value = "5"/>
                <xsd:pattern value = "[a-zA-Z]{2}\-[a-zA-Z]{2}"/>
              </xsd:restriction>
            </xsd:simpleType>
          </xsd:attribute>
        </xsd:complexType>
      </xsd:element>
    </xsd:sequence>
    <xsd:attribute name = "name" use = "required" type = "xsd:string"/>
  </xsd:complexType>
</xsd:element>
</xsd:sequence>
<xsd:attribute name = "search_attr_name" use = "required" type = "xsd:string"/>
<xsd:attribute name = "search_attr_type" use = "required">
  <xsd:simpleType>
    <xsd:restriction base = "xsd:string">

```

```
        <xsd:enumeration value = "string"/>
        <xsd:enumeration value = "number"/>
        <xsd:enumeration value = "date"/>
    </xsd:restriction>
</xsd:simpleType>
</xsd:attribute>
</xsd:complexType>
</xsd:element>
</xsd:sequence>
</xsd:complexType>
</xsd:element>
</xsd:schema>
```

Altering the Crawler Java Classpath

The Ultra Search crawler is a pure Java application that runs in a Java virtual machine. A Java virtual machine uses the Java classpath to find classes during runtime. When Ultra Search is installed, the default crawler classpath is stored in the database. Whenever a new Ultra Search instance is created, this default classpath is copied and used as the crawler classpath for that specific instance.

Reasons for Altering the Crawler Java Classpath

Usually, you do not need to alter the crawler Java classpath. However, there are certain reasons for you to do so. One reason could be to replace the JavaMail reference implementation with a third party JavaMail implementation.

Difference Between the Crawler Classpath and the Remote Crawler Classpath

The crawler classpath is the classpath of a crawler that runs on the same host as the Ultra Search backend (server component). However, Ultra Search allows remote crawlers to be run on other hosts for scalability.

Remote crawler activation uses Java remote method invocation (RMI) technology. As a result, the classpath setting of a remote crawler is inherited from the classpath settings of the RMI registry and RMI daemon.

See Also: "Using the Remote Crawler" on page 4-6

Altering the Crawler Java Classpath on the Ultra Search Server Host

1. Log on to the host where the Ultra Search backend (server component) is installed. Locate the file `$ORACLE_HOME/ultrasearch/admin/wk0addcpath.sql`.
2. Using SQL*Plus, run the `wk0addcpath.sql` script as the `WKSYS` super-user or as a database user that has been granted the super-user privileges. (This script only updates the `CRAWLER_CONFIG_DEFAULT` table. You also need to reconfigure your crawlers to get the `WK$CRAWLER_CONFIG` table updated correctly.)
3. When prompted, specify whether you want to alter the default classpath or an instance-specific classpath. Altering the default classpath causes all subsequently created instances to use that classpath. Existing instances are not modified.
4. When prompted, enter the Ultra Search instance name if you are attempting to modify an instance-specific classpath. If you are modifying the default classpath, then you do not need enter anything here.
5. When prompted, specify whether you want to update the entire classpath or append to it. Appending to a classpath adds entries to the beginning of it. Usually, earlier entries in the classpath override later entries in the case of duplicate classes.
6. When prompted, enter the new classpath if updating the entire classpath. If you are appending one or more directories or library files to the classpath, then enter these separated by the classpath separator for the platform where the Ultra Search backend is installed (the colon on UNIX platforms, or the semicolon on Windows).

Altering the Crawler Java Classpath on a Remote Crawler Host

1. Log on to the remote crawler host where the Ultra Search middle tier is installed. On a UNIX computer, locate and open the file `$ORACLE_HOME/ultrasearch/tools/remotecrawler/scripts/unix/define_env`. On a Windows computer, locate and open the file `$ORACLE_HOME/ultrasearch/tools/remotecrawler/scripts/winnt/define_env.bat`.
2. The `define_env` file specifies all environment settings used by the RMI subsystem. To alter the classpath, use a text editor to modify the `APPLICATION_CLASSPATH` variable.

3. Restart the RMI subsystem for these changes to take effect.

See Also: "Using the Remote Crawler" on page 4-6 for more details on starting up the RMI subsystem

Customizing the Query Syntax Expansion

9.0.1

Oracle Ultra Search uses the Oracle Text engine to index and search documents. When an end user specifies a certain query string, Oracle Ultra Search takes that string and transforms it into an Oracle Text query expression. This process is called query syntax expansion.

Oracle Ultra Search provides a default query syntax expansion implementation. The code for this implementation is contained in the `WK_QUERYEXP` PL/SQL package. It can be viewed in the `$ORACLE_HOME/ultrasearch/admin/wk0queryexp.pkb` file on the Oracle Server host.

This appendix describes how to customize the query syntax expansion implementation to suit your organization's preferences.

Default Query Syntax Expansion Implementation

The default query syntax expansion implementation directly affects the following.

- The way the end user enters a query string (known as the *end user query syntax*)
- The way the documents matching the query are scored (known as *scoring*)
- The way the end user's query string is transformed into an Oracle Text query string (known as the *expansion rules*)

End User Query Syntax

The end user query syntax defined by the default query syntax expansion implementation is similar to the standard text query syntax employed by most search engines on the Web.

- **Token:** A token is a string enclosed in double-quotes ("). It can be a single word or a phrase.
- **Operators:** The default implementation defines three operators. They are the [+], [-], and [*] operators. These operators are defined by the default implementation. Change these operators to whatever you prefer in your own custom implementation.

The plus operator [+] specifies that the token immediately following it must appear in all documents included in the search result.

The minus operator [-] specifies that the token immediately following it cannot appear in any document included in the search result.

The asterisk [*] operator specifies a wildcard search. It matches zero or more characters. A token starting with the asterisk is ignored. The asterisk can only be specified at the end (right side) or middle of a token. For example, "hel*o" and "hell*" use the asterisk correctly, but "*ello" is unacceptable.

Summary of Rules

The following table summarizes the rules for the Ultra Search end user query syntax:

Note: All end-user query strings are encased in square braces. For example, the end user query string Oracle Applications is notated as [Oracle Applications].

Rule	Description
Single word search	<p>Entering one word finds documents that contain that word.</p> <p>For example, searching for [Oracle] finds all documents that contain the word "Oracle" anywhere in that document.</p> <p>Note: Searching for [Oracle] is not equivalent to [Oracle*].</p>
Multiple word search	<p>Entering more than one word finds documents that each contain any of those words in any order.</p> <p>For example, searching for [Oracle Applications] finds documents that contain "Oracle" or "Applications" or "Oracle Applications."</p>

Rule	Description
Compulsory inclusion [+]	<p>Attaching a [+] in front of a word requires that the word be found in all matching documents.</p> <p>For example, searching for [Oracle + Applications] only finds documents that contain the word "Applications." Note: In a multiple word search, you can attach a [+] in front of every token including the very first token.</p>
Compulsory exclusion [-]	<p>Attaching a [-] in front of a word requires that the word must not be found in all matching documents.</p> <p>For example, searching for [Oracle - Applications] only finds documents that do not contain the word "Applications". Note: In a multiple word search, you can attach a [-] in front of every token except the very first token.</p>
Phrase matching ["..."]	<p>Putting quotes around a set of words only finds documents that contain that precise phrase.</p> <p>For example, searching for ["Oracle Applications"] finds only documents that contain the string "Oracle Applications."</p>
Wildcard matching [*]	<p>Attaching a [*] to the right-hand side of a word returns left side partial matches.</p> <p>For example, searching for the string [Ora*] finds documents that contain all words beginning with "Ora," such as "Oracle" and "Orator." You can also insert an asterisk in the middle of a word. For example, searching for the string [A*e] retrieves documents that contain words such as "Apple", "Ate", "Ape", and so on. Wildcard matching requires more computational processing power and is generally slower than other types of queries.</p>

Scoring

There are three ways documents are matched against an end user query string. These three ways are known as scoring "classes." Documents are scored and ranked higher if they satisfy the requirements for a higher class. Within each class, documents are also ranked differently depending on how well they match the conditions of that scoring class.

Class 1 is the most heavily weighted class. The score is derived from the number of occurrences of a precise phrase in a document. A document that has more instances of the precise phrase have a higher score than another document that has fewer occurrences of the precise phrase.

Class 2 is the next more heavily weighted class. In this class, the closer the tokens appear in a document, the higher the score becomes. For example, an end user query string [Oracle Applications Financials] can result in three documents found. None of the three documents contain the precise phrase "Oracle Applications Financials." However, document X contains the all three tokens "Oracle", "Applications", and "Financials" in the same sentence separated by other words. Document Y contains the individual tokens in the same paragraph but in different sentences. Document Z contains the same three tokens, but each token resides in different paragraphs. In this scenario, document X has the highest score, because the tokens are closest together. Likewise, Y has a higher score than Z.

Class 3 is the least weighted class. A document that has more tokens gets a higher score. For example, an end user query string [Oracle Applications Financials] can result in three documents found. Document X might contain all three tokens. Document Y might contain the tokens "Oracle" and "Applications" only. Document Z might contain only the token "Oracle." In this scenario, document X has a higher score than Y. Likewise, Y has a higher score than Z.

Expansion Rules

As mentioned earlier, the end user query is expanded to an Oracle Text query. The expanded query string rules are captured in BNF (Backus Naur Form) notation. Again, these rules are the rules that Ultra Search uses as a default query syntax expansion implementation.

The rules that define an expanded query:

```
<expanded query> ::= (<expression> within <title section>)*2, <expression>
<expression> ::= <generic query expression> | <simple query expression>
<generic query expression> ::= ([[ <plus expression>*100 & ]] (<main
expression>)) [ <minus expression> ]
<simple query expression> ::= (<phrase expression>)*2, (<main expression>)
<main expression> ::= (<near expression>)*2, (<accum expression>)
```

The following section contains some terms and their meanings, which explain some of the terms used in the preceding rules:

A <plus expression> is an AND expression of all plus tokens.

A <minus expression> is a NOT expression of all minus tokens.

A <phrase expression> is a PHRASE formed by all tokens in the <main expression>

A <near expression> is a NEAR expression of all tokens but minus tokens.

An <accum expression> is an ACCUMULATE expression of all tokens but minus tokens.

A <simple query expression> is used only when the end user query has multiple tokens and does not have any operator or a double quote. Otherwise, a <generic query expression> is used.

If there is no token that is neither plus token or minus token, then the <plus expression> and the <accum expression> are eliminated.

Examples of Applying the Rules

The following table illustrates how the default query syntax expansion implementation converts end user query strings to Oracle Text compatible query strings.

End User Query String	Expanded Query String Understandable by Oracle Text
[Oracle]	((({Oracle}) within TITLE__31)*2, ({Oracle}))
[Oracle + Applications]	(((((Applications))*10)*10&((Oracle);{Applications})*2, ({Oracle}, {Applications}))) within TITLE__31)*2, (((Applications))*10)*10&((Oracle);{Applications})*2, ((Oracle), {Applications})))
[Oracle - Applications]	((({Oracle})~{Applications}) within TITLE__31)*2, ((Oracle)~{Applications}))
["Oracle Applications"]	((({Oracle Applications}) within TITLE__31)*2, ({Oracle Applications}))
[Ora*]	((({Ora%}) within TITLE__31)*2, ((Ora%)))
[Oracle Applications]	((({Oracle Applications})*2, ((Oracle);{Applications})*2, ({Oracle}, {Applications}))) within TITLE__31)*2, ((Oracle Applications)*2, ((Oracle);{Applications})*2, ({Oracle}, {Applications})))

Customizing the Rules

Customize this expansion to suit your organization's purposes by defining and implementing your own query syntax expansion. To do so, you should understand the requirements of Oracle Text queries. The details of Oracle Text queries are beyond the scope of this document.

See Also:

- *Oracle Text Application Developer's Guide*
- *Oracle Text Reference*

To customize Ultra Search to use your own implementation of the query syntax expansion, modify the `WK_QUERYEXP` package. In that package are two PL/SQL functions that you must edit. The functions are `expand_main` and `expand_attr`. `Expand_main` is applied to the query string entered by the end user. `Expand_attr` is applied to each search attribute specified in an advanced search. The return value of each `expand_attr` function is appended to the return value of the `expand_main` function. This resultant query string is what's given to Oracle Text to query on.

The `expand_main` Function

This takes the query string entered in the basic search box or advanced search box and converts it to an Oracle Text query string according to your custom query syntax expansion implementation rules.

```
CREATE OR REPLACE FUNCTION expand_main(query varchar2)
RETURN varchar2
AS
    newqry varchar2(4000);
BEGIN
    newqry := <Convert the input query string into an
              Oracle Text query string according to
              your custom rules>
    return newqry;
END;
```

The `expand_attr` Function

This is applied to each search attribute in an advanced search. It takes each attribute and converts it to an Oracle Text query string according to your custom query syntax expansion implementation rules.

```
CREATE OR REPLACE FUNCTION expand_attr(query varchar2)
RETURN varchar2
AS
    newqry varchar2(4000);
BEGIN
    newqry := <Convert a search attribute into an
              Oracle Text query string according to
```



```

        your custom rules>
return newqry;
END;

```

Note: All customized functions are instance-specific and should be defined in the schema of the Ultra Search instance user.

Functions should be executed with definers-rights.

Example of Combining Values

The following example illustrates how the default query syntax expansion implementation converts the end user query string Oracle Applications to an Oracle Text compatible query string. The additional clause added by the introduction of the two search attributes is highlighted in bold.

End User Query String	Expanded Query String Understandable by Oracle Text
[Oracle Applications]	((({Oracle Applications}) *2, (({Oracle}; {Applications}) *2, ({Oracle}, {Applications}))) within TITLE__31) *2, (({Oracle Applications}) *2, (({Oracle}; {Applications}) *2, ({Oracle}, {Applications}))))
[Oracle Applications] with the Title attribute restricted to "MyTitle" and the Author attribute restricted to "MyAuthor"	((({Oracle Applications}) *2, (({Oracle}; {Applications}) *2, ({Oracle}, {Applications}))) within TITLE__31) *2, (({Oracle Applications}) *2, (({Oracle}; {Applications}) *2, ({Oracle}, {Applications})))) & ((({MyTitle}) WITHIN TITLE__31) & ((({MyAuthor}) WITHIN AUTHOR__32)) *10) *10

Index

A

access control lists, 1-6
access URL, 6-3, 8-20, 8-30
administration groups, 1-11
authentication, 1-12
 single sign-on, 1-12, 5-3, 7-3

C

caching documents, 6-5
crawler, 6-2
 classpath, B-1
 crawler agents, 6-3
 crawling process, 6-3
 data sources, 6-2
 overview, 6-2
 parameters, 7-2, 7-42
 remote crawler, 7-16
 settings, 6-2, 7-12
 statistics, 7-16
crawler agent
 API, 8-19
 functionality, 8-21
 sample agent files, 8-24
 setting up, 8-24
 smart agent, 8-21
 standard agent, 8-20
crawler agent API, 1-3
crawler agents, 1-5
CTXSYS user, 3-4

D

data groups, 7-3, 7-38
data harvesting mode, 1-11
data sources, 7-20
 email, 7-26
 file, 7-27
 synchronizing, 6-3
 table, 7-23
 synchronization, 4-13
 user-defined, 6-3, 7-31
 Web, 7-20
data-sources.xml file, 2-22
DB_CACHE_SIZE parameter, 4-3
DBMS_JOB package, 1-3
default instance, 5-3
display URL, 6-3, 7-22, 7-24, 7-25, 7-28, 8-20, 8-30
document attributes, 1-9, 6-3
domain rules, 8-30

E

email API, 1-4, 8-26, 8-27
Enterprise Manager, 2-14, 2-25, 3-2, 7-3

F

federated search, 1-6
Federator searchlet, 7-30
federator.rar, 7-30

H

HTTPS, 5-3, 7-21, 7-28

I

index
 altering, 3-7, 6-2
 optimizing, 7-37
indexing documents, 6-7
INSO filter, 2-5
instance snapshot, 1-11

J

Java classpath, B-1
JAZN, 5-8
jazn-data.xml, 5-9
JDBC, 2-21, 2-22, 2-23, 2-28, 4-5, 4-6, 4-12, 4-13, 5-8,
 7-15, 8-2, 8-11, 8-12, 8-24, 8-33, 8-34, A-1, A-2
JOB_QUEUE_PROCESSES initialization
 parameter, 3-4

L

list of values (LOV), 1-8, 1-10, 6-4, 7-18, 7-19, 7-43,
 7-45, 8-10, 8-24, A-1

M

metadata, 6-3
 loading, A-1
metadata loader, 1-10
migration logs, 3-14

O

OC4J, 2-12, 2-14, 2-15, 2-16, 2-18, 4-5, 5-8, 8-33
Oracle Internet Directory, 1-11, 2-6, 5-3, 5-8, 7-3
Oracle Text, 1-2, 1-3, 2-2, 2-9, 3-4, 4-10, 6-2, 6-5, 6-7,
 8-3, C-1
OracleAS Portal, 1-5

P

path rules, 7-27, 8-30
PROCESSES initialization parameter, 3-4
proxy server, 7-17

Q

query API, 1-3, 1-8, 8-2
query statistics, 7-40
query syntax expansion, 1-13, 8-3, C-1
query tag library, 8-9
queuing documents, 6-4

R

redo log files
 sizing, 3-2
relevancy boosting, 1-10, 7-39
 limitations, 1-10
remote crawler, 4-6, 6-9
 profiles, 7-16
remote crawler hosts
 installing, 2-27
resource adapters, 1-6
robots exclusion, 1-9, 7-22
robots.txt file, 1-9, 7-22, 8-29

S

sample query applications, 1-7, 8-33
schedules
 data synchronization, 7-33
 index optimization, 7-37
search attributes, 1-9, 7-18
searchlets, 1-6
secure search, 1-6, 2-8
secure searching, 2-6
Single Sign-On Server, 1-5
SORT_AREA_RETAINED_SIZE initialization
 parameter, 3-4
SORT_AREA_SIZE initialization parameter, 3-4
stoplists, 3-7
 default, 3-7
 modifying, 3-8

T

triggers, 4-15

U

Ultra Search

- administration tool, 7-1
 - administrative privileges, 7-43
 - APIs, 1-3
 - backend (server component), 2-4
 - components, 1-2
 - configuration, 1-13
 - configuring, 3-2
 - crawler, 1-2, 6-2
 - default instance, 5-3
 - globalization, 7-43
 - instance
 - default, 2-11
 - instance administrators, 1-12, 5-3
 - instances, 7-6, 7-10
 - creating, 7-7
 - snapshot, 7-7
 - integration with OID, 1-11, 5-8
 - integration with Oracle Application Server, 1-5
 - languages, 7-13, 7-43
 - logging on, 7-3
 - managing users, 7-42
 - metadata loader, A-1
 - middle tier, 2-12
 - on Real Application Clusters, 4-10
 - overview, 1-2
 - remote crawler, 6-9
 - search portlet, 1-8
 - server component, 1-3
 - snapshot instances, 7-8
 - super-users, 1-12, 5-3
 - system requirements, 2-2
 - tuning, 4-2
 - upgrading, 3-11
 - approaches, 3-12
 - users, 7-42
- Ultra Search searchlet, 7-29
- ultrasearch.rar, 7-29
- undo space
 - sizing, 3-3
- URL link filtering, 8-29
- URL link rewriting, 8-30
- URL looping, 4-2

- URL rewriter, 1-9, 7-22, 8-29
 - creating, 8-32
 - using, 8-32
- URL rewriter API, 1-4
- URL submissions, 7-38
- UrlRewriter, 7-22, 8-29

W

- Web crawling, 8-29
- WK_INST default instance, 5-3
- WK_TEST instance administrator, 5-3
- wk0migrate.sql script, 3-13, 3-14
- wk0pref.sql file, 3-7, 6-2
- wk0upgrade.sql script, 3-12, 3-14
- WKSYS database user, 2-5, 2-28, 3-6, 3-8, 3-14, 5-4, 7-4, B-2
 - changing password, 3-2
- WKSYS.WK_QRY package, 4-5
- WKUSER role, 3-6, 7-43

X

- XML DB, 1-7, 2-6

