

# Oracle Ultra Search

Unlock Your Information Assets

*An Oracle Technical White Paper*  
*March 2001*

EXECUTIVE SUMMARY .....	3
1. INTRODUCTION.....	3
2. ULTRA SEARCH FOR ENTERPRISE INFORMATION AND CONTENT.....	4
3. ULTRA SEARCH ARCHITECTURE.....	6
3.1.1 The Ultra Search Crawler Component.....	9
3.1.2 The Ultra Search Query API and Query Application.....	10
3.1.3 The Ultra Search Administration Component .....	12
3.1.4 The Java Email API.....	12
3.2 Ultra Search Methodology .....	13
3.2.1 The Gather Step.....	13
3.2.2 The Analyze Step .....	14
3.2.3 Making Crawling Results Searchable .....	15
3.2.4 The Maintain Step .....	16
5. SUMMARY .....	17

**The proliferation of information has caused chaos inside firewalls, and the resulting difficulty in locating information is causing inefficiencies, and expense.**

## EXECUTIVE SUMMARY

Oracle Ultra Search is an out-of-the-box search solution that provides search across multiple repositories - Oracle databases, IMAP mail servers, HTML documents served up by a Web server, files on disk and many more. Ultra Search enables a 'Portal' search across the content assets of a corporation, bringing to bear Oracles core capabilities of platform scalability and reliability.

This paper gives IT managers and architects an understanding of Ultra Search - it's architecture, it's main features, it's interfaces and the configurations available for it's deployment.

## 1. INTRODUCTION

In the age of the Internet, proliferation of information is causing a new information management crisis for enterprises. Using the World Wide Web, workers become their own information retrieval experts. But searching for the right answers can be more than frustrating:

- Studies predict that by 2006 the amount of information flowing over corporate Intranets will be 200 times what it was in 1998.
- This information is ultimately stored in corporate databases, Web pages, files in various popular document formats and in email or groupware systems. The information that businesses produce, store and use for decision making is scattered across billions of documents and data fragments that reside on many different, and often incompatible, IT servers and systems. Servers are located throughout the country and across the globe.
- Corporate information is distributed across enterprises in both structured and unstructured form - structured relational databases, unstructured Word-processing documents, spreadsheets, presentations.
- As applications demand transactional consistency, coordinated multi-user access, administration and maintenance for content, a natural gradient is created to move more and more information into databases. However, even when multiple databases are involved, searching across databases needs a robust solution.

According to some estimates, the lost time in searching costs companies billions of dollars in lost productivity each year. Bad search can also drive customers to a

competitor's Web sites. A company can have great products and terrific looking Web sites, but if customers and employees can't find the information they are looking for, they have essentially wasted time and money spent on development and promotion.

Oracle Ultra Search solves the problem of finding relevant information across your company's many disparate repositories of information. Ultra Search is an out-of-the-box application built on Oracle's proven 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 many more. Ultra Search uses a 'crawler' to index documents; 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 thus enables a 'Portal' search across the content assets of a corporation without the need for rearchitecting IT topologies, compromising security, or programming against hard-to-use API's.

## 2. ULTRA SEARCH FOR ENTERPRISE INFORMATION AND CONTENT

**Oracle Ultra Search is a multi-repository search solution that leverages the award-winning search quality of Oracle Text.**

Oracle Ultra Search is an out-of-the-box search solution that:

1. Searches text across multiple repositories – Multiple databases, HTML Web pages, Files, IMAP mail servers - and organizes and categorizes the content.
2. Provides the best relevance ranking and globalisation support in the industry.
3. Provides value added Portal functionality – crawling, fielded search and metadata extraction.
4. Presents a Web-style interface where users can specify, for example, to indicate "Oracle +Location -France" to indicate they want to retrieve all documents, database records or email containing the terms Oracle and Location, but not the term France.

Enterprises can benefit by using Ultra Search in many different types of applications:

1. **Portal Search** – Ultra Search offers the most powerful search for Enterprise Portals developed with the Oracle Enterprise Portal Framework. Oracle9iAS Portal customers can use Ultra Search through a 'Portlet' (a portlet is a contained area of Portal page that can be rendered in HTML or any other browser-capable technology). The Ultra Search Portlet provides crawling and universal search over all Ultra Search-supported repositories, including the ability to search the 9iAS Portal repository.

For organizations who want to build their own portal from scratch, Ultra

Search provides a canned, end-user-oriented, web-style search over various corporate databases, HTML pages, IMAP email servers, or filesystem documents. You can either use our 'default' user-interfaces as supplied, or 'embed' Ultra Search in your portal, customizing the look-and-feel to your requirements. Ultra-search also allows you to customize metadata according to the different repositories, and search according to different metadata elements from different repositories.

2. **Web Search for Oracle Text** – Ultra Search is an application built on Oracle Text, Oracle's industry leading text retrieval engine. It provides Oracle Text customers with Web-style searching capabilities without the need for any low-level SQL programming. A significant amount of expertise has gone into translating and tuning web-style queries into underlying SQL-based Oracle Text queries. Ultra Search helps Oracle Text users start that much ahead, for example database applications needing a simple Text search component will find Ultra Search admirably integrated with the Oracle platform.
3. **Library or Archive Search** -- Many organizations with digital libraries, information warehouses or centralized repositories are seeking to convert custom search applications over such repositories to more general, web-based ones. A Library search differs from a Portal search in that the latter seeks a simple search over many dynamically changing sources, whereas the former needs more advanced search over a fewer number of relatively well-defined sources. Ultra Search provides such lower-level API and linguistic access to meet the needs of advanced knowledge workers.
4. **Content Management Platform Search** –. Media organizations creating or publishing content in a collaborative manner need to search across content (Web pages, documents) as it moves through multiple repositories in different stages of the content-management life cycle: from the desktop file of the author to the staged version in a database. Use Ultra Search to build a better search and retrieval system for your documents by integrating Ultra Search with your company's own collaborative content management of document management process. Ultra Search provides both full-text and fielded text retrieval to create a set of indexes tuned for keeping track of your content.

Search can be improved if it can be narrowed down what part of a 'document' a piece of information occurs in - the title, the body, the name of the author and so on. For example, search results for 'London' differ when you look for an author name, versus a title. Generally, different repositories have different such 'metadata' attributes that may be attractive for searching against - databases identify columns and email servers know header/body/attachment.

A flexible metadata mapping methodology that lets you unify diverse repositories in common logical terms for search purposes is one of the big value-adds of Ultra Search. In order to display a uniform set of results ranked by overall relevance,

Ultra Search allows customers to normalize or map the various metadata attributes from various repositories.

The screen shot below shows an example of querying over multiple repositories: A corporate email archive and a database (labeled 'Server Technologies'). The query is narrowed by metadata fields that have been defined to map against both repositories: 'Title' maps against the subject line in emails and against a fielded text column in the database; 'Author' maps against the sender of an email.



**Figure 1: Oracle Ultra Search Query Example.**

The next section takes a look at Ultra Search architecture, followed by detailed looks at the important aspects of its functionality.

### 3. ULTRA SEARCH ARCHITECTURE

Ultra Search integrates proven technologies in the Oracle platform in a simple yet robust architecture. Ultra Search is entirely an Oracle Text application, built using the same public interfaces that are available to users of Oracle Text, but enhanced with considerable expertise in aggregating information for indexing, translating queries for the best quality search, and optimizing operations for scalability.

Oracle Text, in its turn, builds on the Oracle platform using public interfaces such as SQL and the Oracle Extensibility Architecture ODCI interfaces. Few search engines can search databases effectively, handicapping them for dynamic data. Oracle Text is highly integrated with the Oracle database for best interoperability with dynamic data. One key strength of Ultra Search is its ability to serve search for database-backed web-sites, applications, archives, or content-bases located in a single database or spread across multiple databases.

Ultra Search is a client program to the Oracle server at run time. It can be deployed in two configurations – in the server tier or the mid-tier.

**Ultra Search builds on the Oracle platform, using existing and proven public interfaces.**

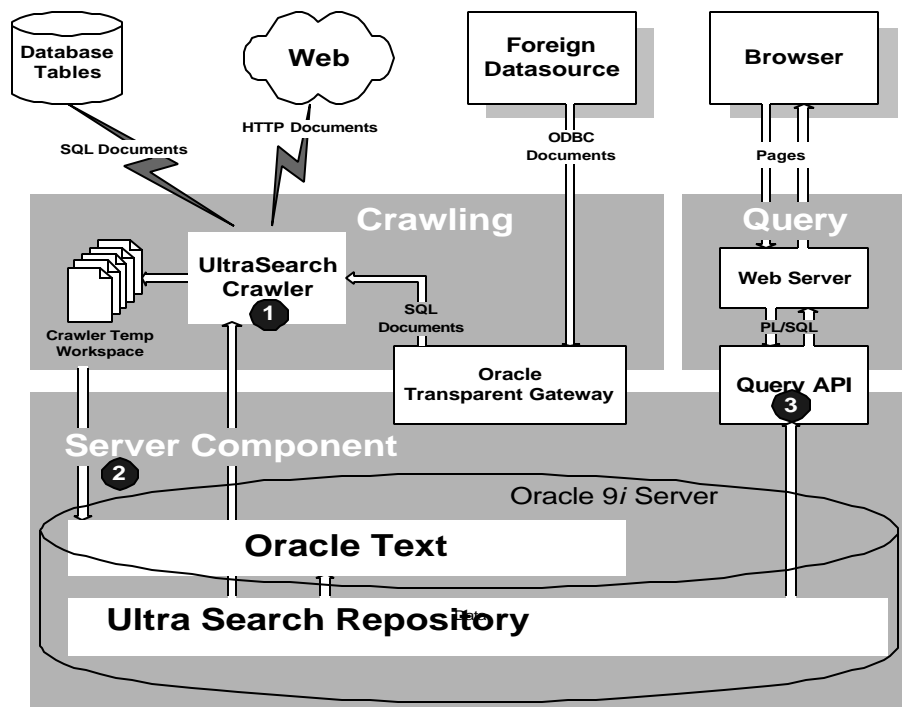


Figure 2: Ultra Search Architecture

Ultra Search is made up of five distinct components:

1. **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 can then be used for querying. The crawler maps link relationships and analyzes them to avoid going in circles and taking wrong turns. 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 the Oracle Text filters to automatically detect the document type and to filter and index the document. See section 3 for more details on supported document types and the filtering process.
2. **Ultra Search Server Component** -- The Ultra Search Server Component consists of a Ultra Search repository, and Oracle9i Text. Oracle9i Text provides the text indexing and search capabilities required to index and query data retrieved from your data sources such web sites, files, and database tables. This component is not visible to users ; it operates as a “black box” that indexes information from the Crawler and serves up the query results.
3. **Query API & Query Application** -- Ultra Search provides Java and PL/SQL APIs for querying indexed data. These APIs return data with and without HTML markup. The HTML markup can help you build the following search engine web interfaces: Basic Search Form, Advanced Search

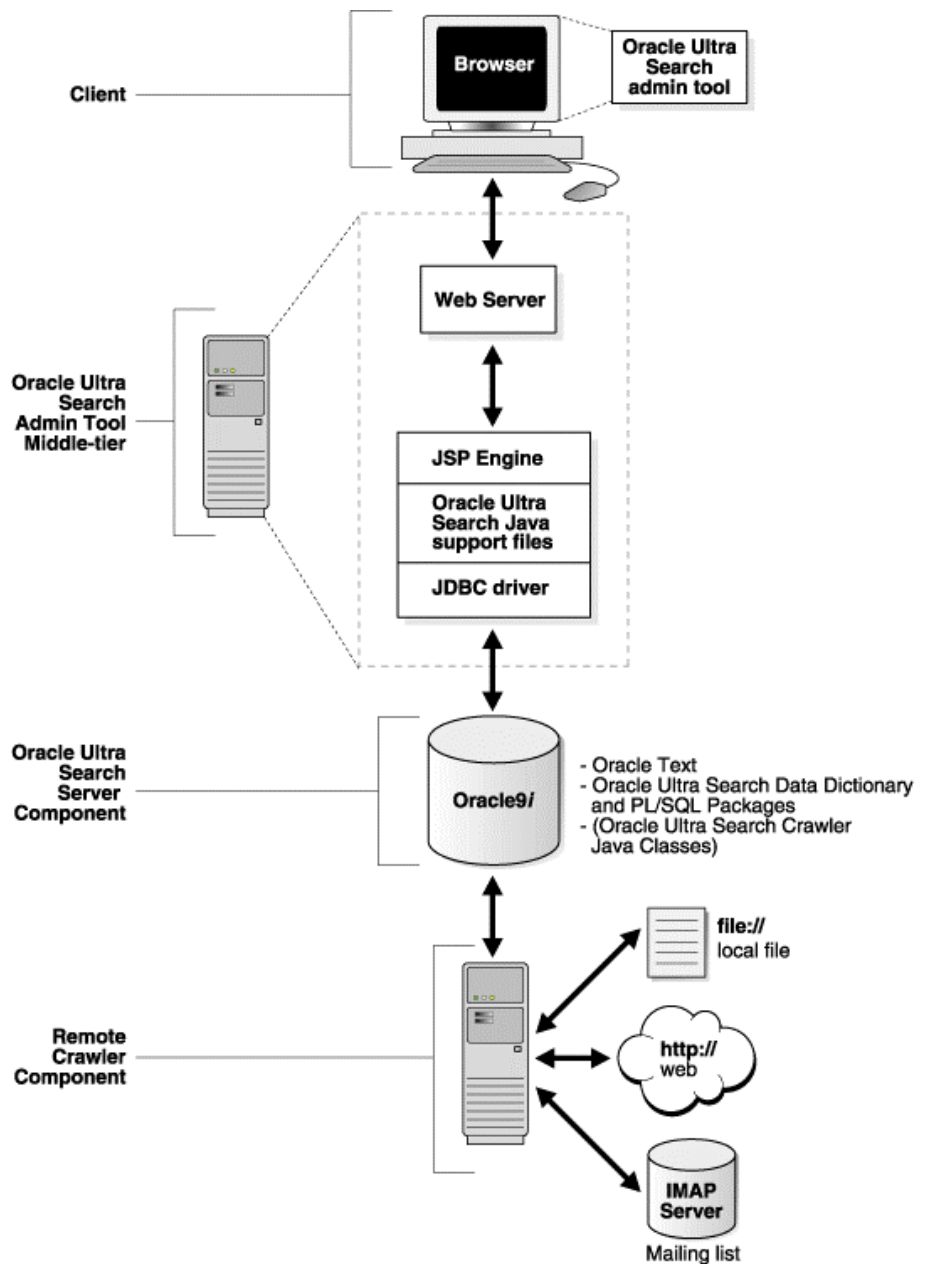
Form, Query Result Display, Help Page, Feedback Page, Register URL. The PL/SQL API requires the Oracle PL/SQL gateway and hence requires Oracle9i iAS. The Java APIs use JDBC connection pooling for scalability.

Ultra Search includes a highly functional query application for users to query and display search results. The query application comes in two versions: a Java Server Page (JSP) version as well as a PL/SQL Server Page (PSP) version. The JSP version will work with any JSP1.0 compliant engine.

4. **Ultra Search Administration Tool and Interface** – The Administration Tool is a Java Server Page (JSP) web application you use to configure and schedule the Ultra Search crawler. The administration tool is typically installed on the same machine as your Web server. You can access the Administration Tool from any browser within your Intranet. The Administration Tool is independent from the Ultra Search Query Application. Therefore, Administration Tool and Query Application can be hosted on different machines to enhance security and scalability.
5. **JAVA Email API** -- Ultra Search provides Java APIs for accessing archived emails. These APIs are used by the Ultra Search Query Application to display emails. These APIs may also be used when building your own custom query application.

The Ultra Search default query interface and the administration tool run in any HTML browser client. The administration tool relies on certain Java classes in the mid-tier. This logical 'mid-tier' can be the same physical machine as the one that runs the database server, or a different one, running Oracle iAS. Finally the Ultra Search database 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 machine or remotely on another machine.

The distribution of Ultra Search components is shown in Figure 3.



**Figure 3: Overview of Ultra Search Components**

### 3.1.1 The Ultra Search Crawler Component

The Oracle Ultra Search crawler is a multi-threaded Java application responsible for gathering documents from the data sources you specify during configuration. The crawler stores the documents in a local file system cache as a temporary workspace during its crawl. Processing the cached data, Oracle Ultra Search creates the index required for querying.

To crawl different repositories, the Ultra Search crawler allows you to define specific 'data sources' (A data source is a logical construct identifying a repository. You can take a single physical repository, such as a database, and map it to multiple data sources. A data source is also the granularity at which you define metadata). Ultra Search knows the following types of data sources:

- **Web Sites** – Define web sites as a data source with the HTTP protocol.
- **Database Tables** – Ultra Search can crawl Oracle databases and other relational databases that support the ODBC standard. Database tables to be crawled can reside in Ultra Search's own database instance or they can be part of a remote, database accessed over a network. To access remote databases, Ultra Search uses 'database links'. Ultra Search allows the crawling of both full text columns and "fielded text" columns. Fielded text columns allow you to map a database column to an Ultra Search attribute (e.g. AUTHOR, TITLE), creating a set of indexes tuned to the content of your database.
- **Files** – Files can be crawled through the file:// protocol. Files must be accessible by each crawler machine either locally or remote over the network. Ultra Search uses the Oracle Text filters to extract text and metadata from documents and automatically identifies document types. See section 3.2.2 for a description of INSO filters and a n overview of supported file types.
- **Emails** --. This feature is useful for crawling mailing lists. Emails sent to a specific email address can be crawled by creating an IMAP email account that subscribes to a mailing list(s). All messages addressed to the email address / mailing list are indexed. Ultra Search can crawl and open email attachments and 'nested' emails such as in email threads.

To maintain fresh, comprehensive search results, Ultra Search uses *synchronization schedules*. Ultra Search lets you gather from multiple Web sites and data sources, each on a separate schedule. Email search results, for example, can be updated continuously, while published content is gathered on a less frequent schedule. Each synchronization schedule can have one or more data sources attached to it.

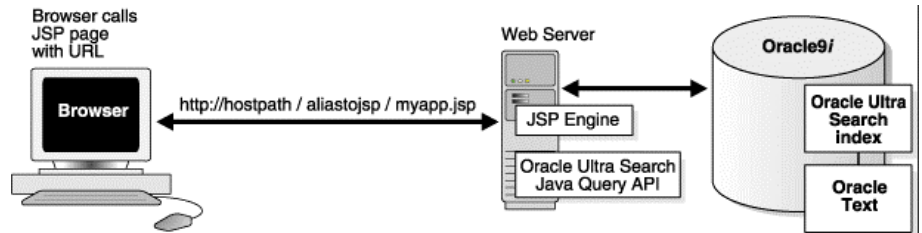
To increase crawling performance, you can set up the Ultra Search crawler to run on one or more machines separate from your database. These machines are called 'remote crawlers'. However, each machine must share cache, log, and mail archive directories with the database machine. To limit the crawling to a specific section of your corporate network or to ensure that crawling does not take wrong turns and follow link relationships that point outside your Intranet, Ultra Search lets you specify so-called 'inclusion' and 'exclusion' domains for crawls.

### 3.1.2 The Ultra Search Query API and Query Application

Oracle Ultra Search provides a flexible, easy-to-integrate query framework by means of a set of query APIs. These APIs can be used from Web applications to retrieve and display query results. Ultra Search APIs are written in both PL/SQL and Java respectively. Therefore, they are compatible with a large spectrum of

web application servers that support either Java Server Pages (JSP version 1.0 and above) or PL/SQL Server Pages (PSP).

Ultra Search Query APIs enable the inserting Ultra Search query input boxes and result lists in any Web application. In addition, these APIs can be used to customize search screens. To include Ultra Search into Java applications, Ultra Search can be called from Java Server Page (JSP) code.

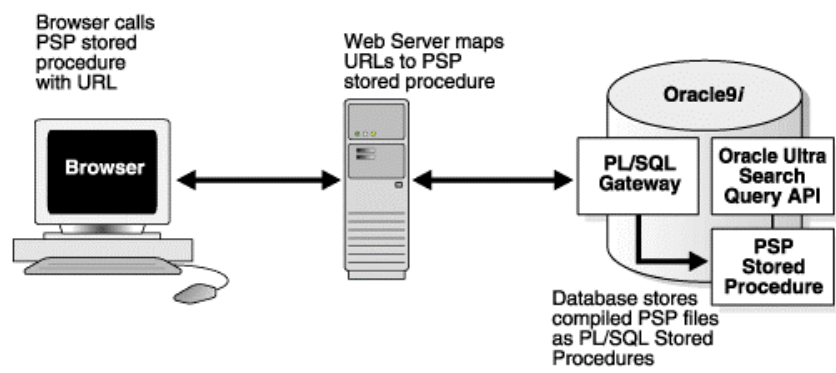


**Figure 4: Illustration of Ultra Search Java Query API**

In this illustration, the browser calls a JSP page with the URL (http://hostpath/...) on the Web server. The Web server, which also contains the JSP Engine and the Ultra Search Java Query API, communicates with Oracle9i, the Ultra Search index and Oracle Text.

To invoke Ultra Search from PL/SQL, procedures in the Ultra Search PL/SQL Query API can be called from a PL/SQL Server Page (PSP). Oracle PSP enables the including of HTML mark up in PL/SQL. PSP files are compiled into an Oracle PL/SQL stored procedure. To run the stored procedure, its URL can simply be typed into a web browser. A PL/SQL gateway (and therefore Oracle9I Internet Application Server) is required in order to map URLs to PSP stored procedures in the database.

Figure 4 below shows how a Web application calls the stored procedure. The Web server maps URLs to the compiled stored procedures. You can then invoke the stored procedures from a browser.



**Figure 5: Illustration of using Ultra Search Query from PL/SQL Server Pages.**

Ultra Search includes a canned, fully functional query application for users to query and display search results. The query application comes in both Java (JSP) and PL/SQL (PSP) versions. The Ultra Search JSP version of the query application also incorporates an email browser for reading and browsing emails.

The Oracle Ultra Search query API is geared towards easy integration with your applications. It provides functionality for

1. Customizing the look and feel of your search per your taste - Ultra Search's Query API provides functions for the presentation of HTML code that can be embedded in your Web application.
2. Retrieving data from the Ultra Search instance, including available data sources, languages and metadata attributes (e.g. AUTHOR, TITLE).

To shorten development cycles, it also includes functionality for encapsulating commonly performed Web development tasks.

### 3.1.3 The Ultra Search Administration Component

The Ultra Search Administration Tool is a web application that allows for:

- **Define Ultra Search instances** – Each Ultra Search instance is identified by name and has its own crawling schedules and index. As many instances as necessary can be created.
- **Manage administrative users** – Ultra Search users can be assigned to manage an instance. Language preferences can be defined for each user.
- **Define crawler parameters** -- Configure and schedule the Ultra Search Crawler
- **Set query options** -- Query options allow users to limit their searches. Searches can be limited to document attributes (e.g. TITLE, AUTHOR) and 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. Each data group consists of one or more data sources.

### 3.1.4 The Java Email API

Oracle Ultra Search enables the retrieving and indexing of emails residing on a server that supports the IMAP4 protocol. Ultra Search defines the concept of an email source, which derives its content from emails sent to a specific email address. When the Ultra Search crawler crawls an email source, it collects all emails that have the specific email address in any of the "To:" or "Cc:" email header fields.

In portal or content search scenarios, the items of interest are usually found in corporate-wide mailing lists and aliases. For example, all customer support issues for Oracle Corp. may be mailed to [support@Oracle.com](mailto:support@Oracle.com). In Oracle Ultra Search,

you create multiple mail sources, where each mail source represents a public email list to which all searchers are assumed to have access.

Ultra Search email crawling and rendering is built on top of the JavaMail API. This enables Ultra Search to provide a Java API for accessing indexed emails. This API enables the retrieving of information such as:

- email header information
- email body content
- and attachments of an email.

The Ultra Search email API allows for including browsing functionality into Java Server 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. Since the source code is viewable, it provides an example for building your own customized email browser.

### **3.2 Ultra Search Methodology**

What steps do you need to follow for using Ultra Search ? The Oracle Ultra Search engine follows four logical steps to provide universal search – *gather, analyze, make queryable, and maintain*. These steps are not novel, and are indeed found in most organizations' business process.

**Error! No topic specified.**

#### **Figure 6: Ultra Search Methodology**

##### **3.2.1 The Gather Step**

Gathering refers to information that exists in structured relational databases and in unstructured files, Word processing documents, spreadsheets, presentations, e-mail, news feeds, Adobe Acrobat files, and Web pages. Ultra Search gathers by this information by “crawling” your corporate Intranet and looking through all the information that exists in the various repositories of your company – databases, Web pages, IMAP mail servers and others. During the gathering process, link relationships are analyzed to avoid going in circles and taking wrong turns. As a result, Ultra Search administrators have an easier time keeping search results complete and up-to-date.



**Figure 7: This screen shot shows the configuring of the Ultra Search crawler for information gathering through the administrator utility. The Web Access page allows users to define starting points for the crawler to begin traveling the corporate Web (e.g. www.oracle.com).**

### 3.2.2 The Analyze Step

In the analyze phase Ultra Search looks at the meaning and structure of gathered information. In order for information to be searched, it must be indexed. During the analyze phase, Ultra Search uses the Oracle Text engine to extract both meaning and structure from the gathered information by creating an integrated index, effectively “normalizing” both structured and unstructured data. Oracle Text indexes contains a complete wordlist along with other information.

During indexing, text and metadata are extracted from documents by Oracle Text INSO filters. This filtering technology automatically identifies document type, invokes the correct filter and produces indexable text and data. Several predefined metadata fields are supported, including author, date, and title. INSO filters include filters for most (150+) popular file types including:

- Microsoft Office Suite 95/97/2000.
- Spreadsheet documents, such as Microsoft Excel and Lotus 1-2-3.
- Word processing files such as Microsoft Word and Corel Word Perfect, including a PDF filter to index Acrobat PDF files.
- Presentation graphics: Microsoft PowerPoint, Lotus Freehand.

Unlike some document management systems, Ultra Search gathering and analyzing is non-intrusive. Instead of physically moving documents, information and documents are analyzed but reside in their original location under their own name.

In typical Web search technologies, hundreds of hits are returned. As the number of repositories increase, the ability to rank relevance of documents decreases.

Ultra Search uses the award winning relevance ranking of Oracle Text to ensure that users consistently find the needle in the haystack.



**Figure 8: The Administrator interface allows you to specify the document types that will be analyzed by Ultra Search and filtered through Oracle Text document filters.**

### 3.2.3 Making Crawling Results Searchable

“Make Searchable” is the function of providing access to all the information that has been indexed in a programmatic fashion. Oracle’s Ultra Search provides both JAVA and PL/SQL API’s for this purpose. Passing a search term into these query APIs locates all relevant documents, whether they are stored on Web servers, databases, or in applications. Customers can use Ultra Search APIs to integrate universal search into their own Web pages or applications.



**Figure 9** Screen Shot of Example Query Screen shows a search for “Performance”-related information. Note Ultra Search relevance rankings appear in red.

### 3.2.4 The Maintain Step

Maintaining ensures that search results are updated continuously. Ultra Search lets you gather from multiple Web sites and repositories, each on a different schedule. IMAP messaging servers, for example, can be updated continuously, while published content is gathered on a less frequent schedule. Ultra Search maintains content by providing easy, intuitive utilities that provide Administrators with an easy way to keep up with new content that is added through growth or aquisition.



**Figure 10:** Screenshot of Ultra Search Administration Utility shows the “Schedule” page where maintenance crawling can be configured.

## 5. SUMMARY

Companies need to eliminate the chaos inside their firewalls. No solution provider is more focused than Oracle to solve that problem.

In summary, Oracle Ultra Search allows you to reduce the time spent finding relevant documents on your company's IT systems:

- It crawls, indexes, and makes searchable your corporate Intranet through a canned, web-style search.
- Provides search without the need for coding against hard-to-use low level API. For advanced users, however, APIs are also exposed.
- It organizes and categorizes content from multiple repositories by extracting valuable metadata that can be used in Portal applications.
- It provides effective search by returning more relevant hits - the best relevance ranking in the industry - and finds what you want.

And it provides the best database integration in the industry.



Oracle Ultra Search  
March 2001  
Author: Stefan Buchta

Oracle Corporation  
World Headquarters  
500 Oracle Parkway  
Redwood Shores, CA 94065  
U.S.A.

Worldwide Inquiries:  
Phone: +1.650.506.7000  
Fax: +1.650.506.7200  
[www.oracle.com](http://www.oracle.com)

Oracle Corporation provides the software  
that powers the internet.

Oracle is a registered trademark of Oracle Corporation. Various  
product and service names referenced herein may be trademarks  
of Oracle Corporation. All other product and service names  
mentioned may be trademarks of their respective owners.

Copyright © 2001 Oracle Corporation  
All rights reserved.