

# How to...

## set up a Web Repository and Crawling it for Indexing

ENTERPRISE PORTAL 5.0 / 6.0

PUBLIC

---

---

ASAP “How to...” Paper



Applicable Releases: EP 5.0 SP5 / EP 6.0 SP0

April 2003

<b>1</b>	<b>SCENARIO .....</b>	<b>1</b>
<b>2</b>	<b>THE STEP BY STEP SOLUTION.....</b>	<b>1</b>
<b>2.1</b>	<b>Creating and Configuring a Web Repository Manager.....</b>	<b>2</b>
<b>2.2</b>	<b>Configuring a Crawler Profile.....</b>	<b>11</b>
<b>2.3</b>	<b>Creating an Index and Assigning a Web RM .....</b>	<b>13</b>
<b>2.4</b>	<b>Monitoring Crawling and Indexing .....</b>	<b>16</b>
<b>3</b>	<b>RESULT .....</b>	<b>19</b>

# 1 Scenario

You want to use SAP Enterprise Portal to search for documents that are stored on various Web sites. The Web server can be reachable in the Internet or in an Intranet.

The scenario assumes that you want to make the Internet pages of a news provider accessible through the portal. The Web server of the news provider is to be checked twice a day for new documents.

The documents can be indexed after the crawler has found them on the Web site. The documents are then available to portal users when they carry out a search. If you are using taxonomies, you can also place these documents into categories. Portal users then access the documents through a taxonomy.

## 2 The Step By Step Solution

You need a Web repository manager (Web RM) to access a Web server. The Web server acts as a data source. You can create a separate Web RM for each Web server, or you can group together several Web servers in a Web repository manager.

See [Creating and Configuring a Web Repository Manager \[Page 2\]](#).

Crawlers that search the Web servers for documents are used to provide the contents of linked Web servers in the portal.

In the KM platform, crawlers and crawler profiles are preconfigured for your use. You have to create a new crawler profile if you want to define how many document levels the crawler should pursue on the Web server.

See [Configuring a Crawler Profile \[Page 11\]](#).

You have to create an index for the Web repository before you can search for documents or classify them. In index administration, select the crawler profile and define how often the Web server is to be searched for new, updated, or deleted documents. After the documents have been crawled they are forwarded to the TREX search engine where they are indexed.

See [Creating an Index and Assigning a Web RM \[Page 13\]](#).

You can use the crawler monitor and the TREX monitor to monitor the status of crawling and indexing.

See [Monitoring Crawling and Indexing \[Page 16\]](#).

As soon as the indexing process has been completed, you can search for documents stored in the Web repository from the portal.

See [Result \[Page 19\]](#).

## Notes

If you are using EP 5.0, you find the KM configuration in the *KM Admin* workset.

In EP 6.0, you find the KM configuration under *System Administration* → *System Configuration* → *Knowledge Management* → *Configuration*.

Configuration parameters that are needed for a simple configuration are highlighted in the table. The other parameters are for enhanced configurations.

## 2.1 Creating and Configuring a Web Repository Manager

This section of the How To guide explains how to create a Web repository manager.

You want to be able to call up documents that are located under the Internet address **www.cnn.com/TECH** in the corresponding Web repository.

Carry out the following steps to create the Web repository manager and the corresponding components.

1. Register the Web server in the KM system landscape.
2. Create a Web site.
3. Optional: Define HTML property extractors.
4. Create a cache for the repository manager.
5. Create a Web repository manager and configure it.

### 1. Register the Web Server in the KM System Landscape

Before you create a Web site, you have to define the URL of the Web server in the form of an HTTP system in the KM system landscape.

To create an HTTP system, choose *Content Management* → *Global Services* → *System Landscape Definitions* → *Systems* → *HTTP System*.

#### Parameters of an HTTP System

Parameter	Required	Description
System ID	Yes	ID of the HTTP system. Caution: Do not use spaces to separate elements of this ID.
Description	No	Description of the HTTP system.
Password	No	Specification of the password needed to access the HTTP system.
Server-URL	Yes	Specification of the URL that targets the HTTP system.

User	No	Specification of the user needed to access the HTTP system.
Same User Domain	No	Specifies whether authentication information for SAP Enterprise Portal is forwarded to the remote server.  Activate this parameter if the HTTP system is to operate in the same user domain as SAP Enterprise Portal.  Deactivate this parameter if you want to access external systems (for example, an external Web site).
<b>max. Connections</b>	No	Specifies the maximum number of connections that a repository manager builds with this system.  The exact number of connections depends on the remote server. We recommend that you minimize the default settings.

Example Configuration for an HTTP System

**New "HTTP-System"**

System ID\*

Description

Password

Re-enter the Password

Server-URL\*

User

Same User Domain

max. Connections

For more information on HTTP systems, see the KM administration guide.

## 2. Create a Web Site

To create a Web site, choose *Content Management* → *Repository Managers* → *Web Sites*.

### Web Site Parameters

Parameter	Required	Description
<b>Name</b>	Yes	Name of the Web site. Caution: Do not use spaces to separate elements of this name.
Display Name	No	The display name for the Web site, which becomes the <code>displayname</code> property of the resource that is displayed in the browser.
<b>Start Page</b>	No	Explicitly specified start page that overrides the default start page of the Web site. Enter the relative part of the URL for the server URL (see example).

System ID (Landscape Service)	Yes	Specifies the <i>System ID</i> entered in the definition of the HTTP system.
System Path	No	Optional suffix for the server URL that is entered in the HTTP system.  This makes it possible to use the same system for different locations in the systems URL namespace. See the screenshot below.

Example Configuration for a Web Site

**New "Web Site"**

Name\*

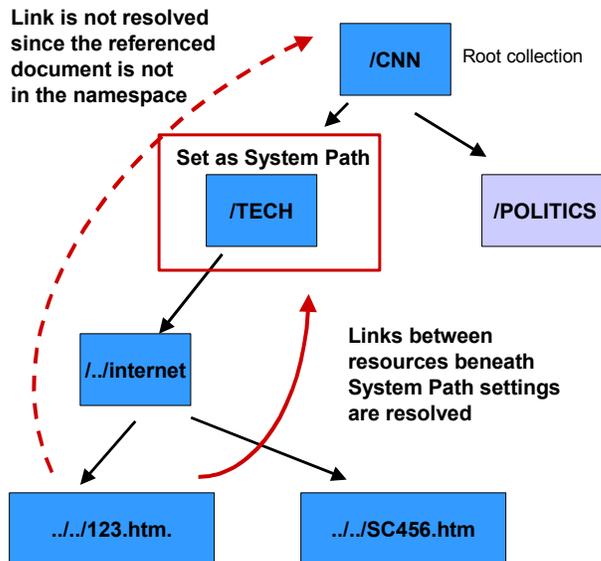
Display name

Start Page

System ID (Landscape Service)\*

System Path

Path Specifications for the Configuration of a Web Site.



If you use the parameter `Start Page` to define the initial point of entry to the repository, there is no limitation on the resolution of links.

If you use the parameter `System Path` (for example, `/TECH`), superordinate links at `/CNN` level are not pursued.

For more information on Web sites, see the KM administration guide.

### 3. Optional: Define HTML Property Extractors

In our example, we want to define certain MIME types that a crawler should not pursue when analyzing the contents of a Web repository. We don't want Shockwave Flash and PDF files to be crawled. Normally, you would also want to exclude ZIP and MP3 files. You enter this information in the HTML property.

To create a new HTML property, choose *Content Management* → *Repository Managers* → *HTML Property Extractors* → *HTML Properties*.

In the parameter `Exclude`, enter the MIME types to be excluded in the form of a *regular expression* (see screenshot). In the parameter `Property Name`, enter the value `embedded-links`.

Example Configuration for an HTML Property

**View "linkfilter"**

Exclude

Namespace

Property Name

Select All

Select First

Select HREF

You now select this HTML property in the definition of an HTML property extractor. To create a new HTML property, choose *Content Management* → *Repository Managers* → *HTML Property Extractors* → *HTML Property Extractor*.

Example Configuration for an HTML Property Extractor

**View "Linkfilter"**

HTML Properties\*

Name
linkfilter

Seite 1/1

You reference the HTML property extractor later on in the configuration for the Web repository manager.

For more information on HTTP property extractors, see the KM administration guide.

#### 4. Create a Cache for the Repository Manager

You create a separate cache for the Web repository manager.

To do this, choose *Content Management* → *Utilities* → *Caches* → *Memory Cache*.

Example Configuration for a Cache

View "ca_CNN"	
Restart lifetime on access	<input type="checkbox"/>
Singleton	<input checked="" type="checkbox"/>
Assumed Entry Size*	<input type="text" value="10000"/> bytes
Capacity*	<input type="text" value="10000"/> entries
Default Time-to-Live (0=infinite)*	<input type="text" value="7200"/> seconds
Max Cache Size (0=unlimited)*	<input type="text" value="10000000"/> bytes
Max Entry Size (0=unlimited)*	<input type="text" value="10000"/> bytes
<input type="button" value="Edit"/> <input type="button" value="Close"/>	

For more information on caches, see the Cache How To and the KM administration guide.

#### 5. Create a Web Repository Manager

To create a Web repository manager, choose *Content Management* → *Repository Managers* → *Web Repository*. In the configuration, you assign the Web site that you just created, the cache, and the HTML property extractor to the Web repository manager.

You have to enter a proxy server and port if

- a) No access is possible to the Web site in the Intranet without them
- b) You want to take advantage of the caching benefits of the proxy

You do not need to make entries for the parameters *Start Page*, *System ID* (Landscape Service) and *System Path*. These parameters are only valid for *simple* Web repositories. If you are using Web sites, you do not use these parameters.

Enter the following parameters in the configuration for a Web repository manager:

### Web Repository Manager Parameters

Parameter	Required	Description
<b>Name</b>	Yes	Name of the repository manager.
Content Cache Directory	No	Specifies a directory in the file system in which the resource contents for the cache are stored.  If no directory is specified, the resource contents are stored in the database.
Description	No	Description of the repository manager
<b>Prefix</b>	Yes	The URI prefix for which the manager is registered. This specification is entered in the list in the root directory.  Note that you must enter the prefix with a forward slash, for example, <code>/WEB REPOSITORY</code> .
<b>Proxy Host</b>	No	Hostname of an HTTP proxy to use for this repository  If the proxy server requires authentication, use the <code>Proxy System ID</code> parameter.
<b>Proxy Port</b>	No	Port number where the proxy is reachable.
Proxy System ID (Landscape Service)	No	Proxy server ID from the system landscape service  You should use this parameter instead of <code>Proxy Host</code> when the proxy server requires authentication.
Start Page	No	Explicitly specified start page that overrides the default start page of the Web site.
System ID (Landscape Service)	No	HTTP system specification that is registered in the CM system landscape. You specify this parameter only if the repository is to map a <b>single Web site</b> . In this case, you must not select any of the Web sites that may be presented as options for the <code>Web Sites</code> property.
System Path	No	Optional suffix for the server URL of the remote system. The server URL is a property of the (landscape) system. <code>System_Path</code> enables you to use the same system for different paths in the URL namespace.
Cache Content Persistently	No	Specifies whether resource content is cached.  Needs a parameter specification in the <code>Connection Pool</code> parameter.  Depending on the <code>Content Cache Directory</code> , the content is stored either in the database or in the file system.

Case-Sensitive URI Handling	No	Determines whether the repository manager takes account of lowercase and uppercase text in resource URLs.  If this is active, the system distinguishes between <i>INDEX.HTM</i> and <i>index.htm</i> , for example.
Dynamic	No	Determines whether the repository can dynamically add new remote HTTP servers to its namespace.
Filenames	No	Determines whether the repository generates resource names that are valid file names.  Links inside an HTML page also use such names. The default value is <i>false</i> .
<b>Send Events</b>	No	Specifies whether the repository sends events when operations such as <i>delete</i> and <i>update content</i> are performed.  The repository sends events if this parameter is activated.
<b>Use System Default Proxy Settings</b>	No	Determines whether the Java Virtual Machine (VM) standard values <code>http.proxyHost</code> and <code>http.proxyPort</code> are taken into account.  You can set these values when starting the VM. If this parameter is activated, the settings apply to all Web repositories provided you have not configured any other settings for <code>Proxy Host</code> or <code>Proxy System ID</code> .
External Server URI Handling	Yes	Determines how URIs (links) inside HTML pages that point to external servers (outside the scope of this web repository) are handled.  <i>none</i> : URIs are transferred unchanged.  <i>rewrite</i> : URIs are rewritten to point to URIs on this server wherever possible.  <i>report</i> : Like <i>rewrite</i> , but the newly written URIs are also transmitted to the property <code>embedded-links</code> for the resource.
Property Search Manager	No	Selection of manager for property search.  It is used by services and applications that need to find resources based on their properties.  <i>Choose Standard Property Search Manager.</i>
Cache Stale Timeout	No	If defined, determines the time in milliseconds after which old resources are deleted from the database.  The lifetime of cached resources is determined by the <code>Cache Timeout</code> parameter.  After the timeout, they are <i>stale</i> . A resource that is older than the specification in the <code>Cache Stale Timeout</code> parameter is removed from the cache.

Cache Timeout	No	Timeout in milliseconds for resources in the cache. The timeout determines the amount of time for which cached resources (content and properties) are not refreshed. The optional value is determined by the update frequency of the remote site. However, the timeout value should be shorter than the update interval of the remote location, so that the cache is updated with the latest content. For resources supplied with an "Expires" header by the remote HTTP server, the timeout is added to the expiry date.
HTTP Timeout	No	Timeout in milliseconds after which operations on the server are aborted. The default value is 180000 ms (that is, three minutes).
<b>Repository Services</b>	No	Identifiers of the repository services you want to use with the repository. The <code>properties</code> repository service needs to be activated so that documents can be classified.
<b>Web Sites</b>	No	Selection of repositories stored in the root directory of the Web repository manager.
ACL Manager Cache	No	Cache identification for resource ACLs. This parameter is required if an ACL security manager is specified in the <code>Security Manager</code> parameter. You can choose the preconfigured cache <code>ca_rsrc_acl</code> .
<b>Connection Pool</b>	No	Identifier of a JDBC connection pool to be used by the Web repository manager for storing resource properties persistently in the database. This enables delta-crawling of the Web repository. You can use a memory cache and connection pool together. Choose <code>dbconrep</code> if you do not want to create your own JDBC connection pool.
<b>HTML Property Extractors</b>	No	Identifier of an HTML property extractor to be used with this repository manager for extracting additional resource properties.
<b>Memory Cache</b>	Yes	Identifier of the cache to be used by the Web repository manager for caching both content and properties of generated resources.
Security Manager	No	Specification of the security manager that controls access to repository contents. If you want the Web repository to perform an authorization check when resources are accessed, you need to specify a security manager. Generally, the <code>AcISecurityManager</code> should be used for Web repositories.

Example Configuration for a Web Repository Manager

**View "Web-RM-How-To"**

Content Cache Directory	<input type="text"/>								
<i>Description</i>	used for the How-To								
<i>Prefix (must start with /)*</i>	/web-repository-how-to								
Proxy Host	proxy								
Proxy Port	8080								
Proxy System ID (Landscape Service)	<input type="text"/>								
Start Page	<input type="text"/>								
System ID (Landscape Service)	<input type="text"/>								
System Path	<input type="text"/>								
<i>Active</i>	<input checked="" type="checkbox"/>								
Cache Content Persistently	<input checked="" type="checkbox"/>								
Case-Sensitive URI Handling	<input checked="" type="checkbox"/>								
Dynamic	<input type="checkbox"/>								
Filenames	<input type="checkbox"/>								
<i>Send Events</i>	<input checked="" type="checkbox"/>								
Use System Default Proxy Settings	<input type="checkbox"/>								
External Server URI Handling*	report								
<i>Property Search Manager</i>	Standard Property Search Manager								
Cache Stale Timeout	<input type="text"/> millisecond								
Cache Timeout	180000 milliseconds								
HTTP Timeout (0=no timeout)	180000 milliseconds								
<i>Repository Services</i>	<table border="1"> <thead> <tr> <th>Name</th> </tr> </thead> <tbody> <tr><td>accessstatistic</td></tr> <tr><td>comment</td></tr> <tr><td>discussion</td></tr> <tr><td>feedback</td></tr> <tr><td>logger</td></tr> <tr><td>personalnote</td></tr> <tr><td>properties</td></tr> </tbody> </table> <p>Seite 1/1</p>	Name	accessstatistic	comment	discussion	feedback	logger	personalnote	properties
Name									
accessstatistic									
comment									
discussion									
feedback									
logger									
personalnote									
properties									
<i>Web Sites</i>	<table border="1"> <thead> <tr> <th>Name</th> </tr> </thead> <tbody> <tr><td>CNN-TECH</td></tr> </tbody> </table> <p>Seite 1/1</p>	Name	CNN-TECH						
Name									
CNN-TECH									
<i>ACL Manager Cache</i>	Not Set								
Connection Pool	dboon_rep								
HTML Property Extractors	Linkfilter								
Memory Cache*	ca_CNN								
<i>Security Manager</i>	Not Set								

For more information on Web repositories, see the KM administration guide.

## 2.2 Configuring a Crawler Profile

This section of the How To guide describes how you create a new crawler profile.

Several crawler profiles are available in the system. However, by default, the number of document levels that are searched is limited (`Depth = -1`). A restriction is recommended in most cases since you otherwise have no control over the size of the Web repository.

To create a new crawler profile, choose *Content Management* → *Global Services* → *Crawler Profiles*.

### Crawler Profile Parameters

Parameter	Required	Description
<b>Name</b>	Yes	Name of the crawler profile.
<b>Description</b>	Yes	Description of the crawler profile. This description is used for the crawler profile listing in index administration.
Case-Sensitive	No	Specifies whether the crawler distinguishes between lowercase and uppercase for resources.  If it is activated, <code>/web/page/index.html</code> and <code>/web/page/Index.HTML</code> are handled differently, for example.
Versions	No	Boolean parameter that determines whether versions of a resource are included in the results set
<b>Hierarchy Mode</b>	Yes	Specifies how hierarchies are crawled. Choose the setting <code>auto</code> for Web repositories.
<b>Depth</b>	Yes	Number of recursion levels applied in the crawling process.  For example, a recursion level of '2' means that starting from a given document, all documents referenced by hyperlinks in the start document and all documents referenced in turn in those documents are included in the results set.  In the case of a crawler that is used for a Web repository, the specification refers to links. The specification should not be less than the values in the parameters <code>External Depth</code> and <code>Internal Depth</code> .  Never choose <code>-1</code> for a dynamic Web repository.
<b>External Depth</b>	Yes	Number of recursion levels used if the crawler finds links to documents in another repository/on another Web server.  If you enter 0, external links are not pursued.  The specification 1 is useful for Web repositories, since in this case links to further information on another server are pursued.
<b>Internal Depth</b>	Yes	Number of recursion levels used if the crawler finds links to documents in the same repository/on the same Web server.  If you enter 0, internal links are not pursued.

Max Content Size	Yes	<p>Specifies the maximum size (in KB) of the content of the resources that register with the TREX components.</p> <p>-1 means no limit. 0 means that only resources are registered, and not their content. This is useful if large documents (for example, pdfs), are not to be indexed.</p>
Priority	Yes	<p>Integer between 1 and 6 that denotes the priority of the crawler thread(s) on the portal server.</p> <p>Priority = 1 (lowest priority) Priority = 6 (highest priority)</p>
Request Pacing	Yes	<p>Integer that determines the size of time lags between successive requests the crawler sends to the Web server whose content is being crawled.</p> <p>The default value is '0', meaning the crawler sends successive requests without any delay.</p>
Time Limit	Yes	<p>The time interval in seconds after which the crawler terminates the crawling process.</p> <p>The default value is -1, meaning that no time limit is set.</p>
<b>Crawler Type</b>	Yes	<p>Specifies the crawler type.</p> <p>Choose <i>web</i> or <i>webdb</i> to crawl a Web repository.</p> <p><i>web</i>: This crawler type is fast but requires a lot of memory.</p> <p><i>webdb</i>: This is slower than <i>web</i> but requires less memory since the crawled resources are stored in the database. Only this crawler type enables delta crawling.</p>

## Example Configuration for a Crawler Profile

**View "Crawler for How-To"**

Description *	Crawler for How-To (Depth=2)	
Case-Sensitive	<input checked="" type="checkbox"/>	
Versions	<input type="checkbox"/>	
Hierarchy Mode *	auto	
Depth (-1=unlimited) *	2	Recursion Levels
External Depth (-1=unlimited) *	0	Recursion Levels
Internal Depth (-1=unlimited) *	2	Recursion Levels
Max Content Size (-1=unlimited) *	-1	KB
Priority *	5	1 (low) to 10 (high)
Request Pacing *	0	Milliseconds
Time Limit (-1=unlimited) *	-1	Seconds
Crawler type *	webdb	

For more information on crawlers, see the KM administration guide.

### 2.3 Creating an Index and Assigning a Web RM

Carry out the following steps to crawl the Web repository for indexing. After the indexing process has been completed, you can start a search on documents that are stored in the Web repository.

The following steps are necessary:

1. Create an index.
2. Assign the Web repository to the index.
3. Select the crawler profile that is to be used to search the repository.
4. Define a plan for indexing the repository.

In EP 5.0, index administration is located in the *KM Admin* workset.

If you are using EP 6.0, it can be found under *System Administration* → *System Configuration* → *Knowledge Management*.

**1. Create an Index**

Choose *Create* in index administration.

Enter the required information and then save your entries.

Field	Description
<b>ID</b>	Enter a unique ID for the index.
<b>Name</b>	Enter a name to be displayed in the search options later on ( <i>Show Indexes</i> ).
<i>Group</i>	An index group includes several indexes. The name specified is displayed in the search options.
<b>Crawler Profile</b>	Select the crawler profile that you just created.
<b>Service</b>	Select an index service. In our example, a pure search index is used. If you want to use the Web repository in a taxonomy, choose a combined search and classification index.

Example Configuration for an Index

**Index Administration**

**how-to-index**

**Properties**

Data Sources

[Back](#)

**Properties**

ID

Name

Group   [Add](#)

Crawler Profile  [Define Schedule](#)

Service

Service Type

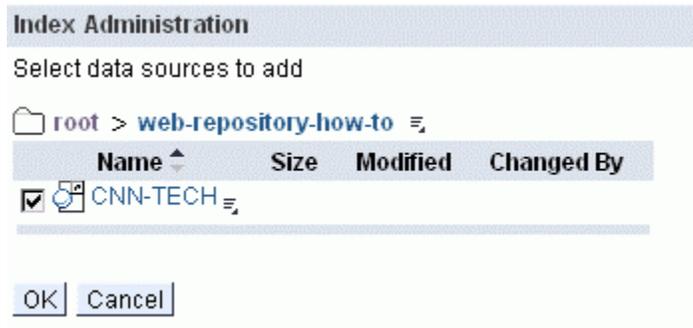
fuzziness

[Save](#)

## 2. Assign a Web Repository/Site to an Index

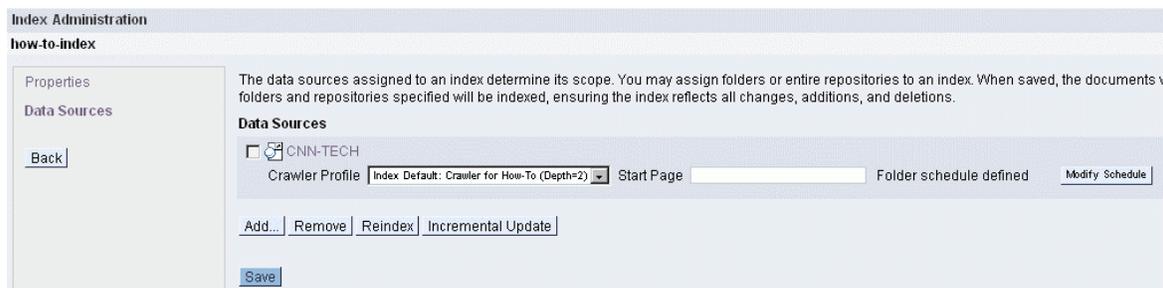
In the left-hand side of the index administration window, choose *Data Sources*. Navigate to the Web repository and select the Web site.

If the Web repository contains multiple Web sites, you can select the Web repository itself instead of the individual Web sites.



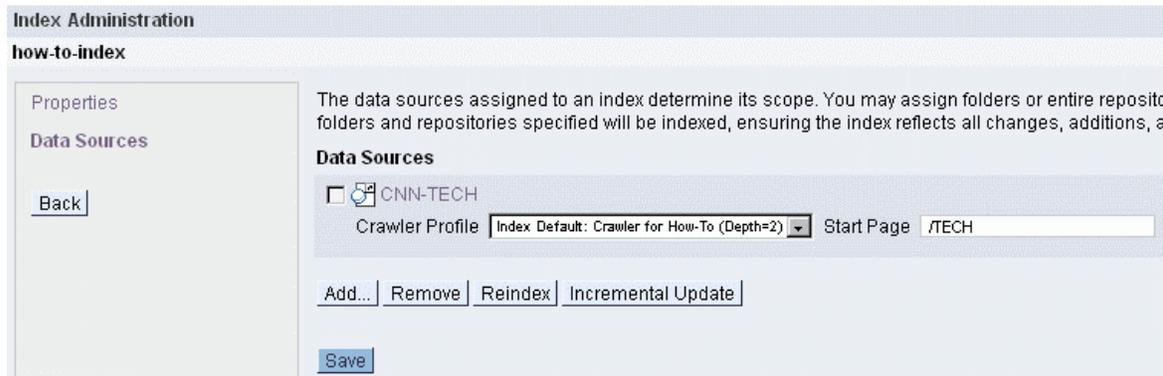
## 3. Select a Crawler Profile

Select the crawler profile that is to be used to search the Web site.



## 4. Specify the Start Page

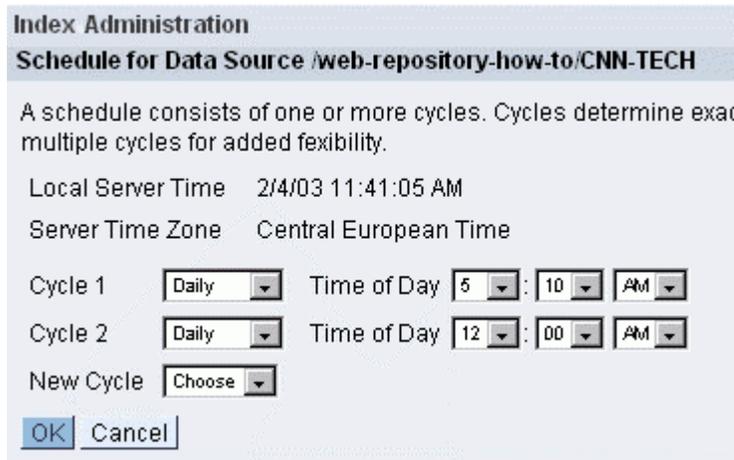
If you entered a start page in the definition of the Web site, you should also enter this start page here.



## 5. Define a Plan for Searching and Indexing the Repository

You can determine the times at which data sources are to be crawled. In our scenario, this takes place twice a day.

Call up the *Define Schedule* function and enter the required time schedule information.



The screenshot shows a dialog box titled "Index Administration" with a subtitle "Schedule for Data Source /web-repository-how-to/CNN-TECH". Below the subtitle, there is a descriptive text: "A schedule consists of one or more cycles. Cycles determine exact multiple cycles for added flexibility." The dialog displays the "Local Server Time" as "2/4/03 11:41:05 AM" and the "Server Time Zone" as "Central European Time". There are two cycle configuration rows. "Cycle 1" has a frequency dropdown set to "Daily", and its "Time of Day" is set to "5:10 AM". "Cycle 2" also has a frequency dropdown set to "Daily", and its "Time of Day" is set to "12:00 AM". A "New Cycle" dropdown is set to "Choose". At the bottom left, there are "OK" and "Cancel" buttons.

For more information on indexes, see the KM administration guide.

## 2.4 Monitoring Crawling and Indexing

The KM platform allows you to monitor the crawling and indexing of Web repositories. You use the crawler monitor and TREX monitor to do this.

The crawler starts shortly after an index is created. You can monitor the status of indexing using the TREX monitor as soon as the crawling process has finished.



Note that the crawling and indexing processes can last several hours. The time required is dependent on the extent of the Web site.

In EP 5.0 the crawler monitor and TREX monitor can be found in the *KM Admin* workset.

If you are using EP 6.0, they can be found under *System Administration* → *Monitoring* → *Knowledge Management*.

### Using the Crawler Monitor

Call up the crawler monitor. You see an overview of the crawlers. This includes the crawler for crawling the Web repository CNN/TECH. The table below contains information on the status of the crawlers.

Status *Running*: The crawler has been started. The Web server is currently being searched.

Status *Completed*: The crawling process has been completed. Documents found are now indexed by TREX.

**Crawler Monitor**  
 Crawlers are processes that retrieve documents and analyze their contents in order to locate additional documents for processing. Several crawler Status information for crawlers that complete is only available for a limited amount of time (30 minutes by default).

Column Group Status Sort By Name ↑↓ Refresh

	Name, Start	Start Path	Status	Start, Stop	Stopped	Scheduled
<input type="checkbox"/>	IndexCrawler: Crawler for How-To-Index	CNN.com - Technology	Running	2/4/03 11:44:05 AM	-	

Slow Down Speed Up Suspend Resume Stop Restart Delete

Last Refreshed at 11:44:36 AM

You can call up more statistics on the crawler. To do this, choose *Statistics* under *Column Group*. An overview shows you the number of documents found.

**Crawler Monitor**  
 Crawlers are processes that retrieve documents and analyze their contents in order to locate additional documents for processing. Several crawlers with the same name may be act Status information for crawlers that complete is only available for a limited amount of time (30 minutes by default).

Column Group Statistics Sort By Name ↑↓ Refresh

	Name, Start	Start Path	Documents	New	Changed	Deleted	Errors	Elapsed Time	Average Time	Last Time	In Progress
<input type="checkbox"/>	IndexCrawler: Crawler for How-To-Index	CNN.com - Technology	53	-	-	-	1	00:01:20.826	00:00:01.525	0.581	551 / 0

Slow Down Speed Up Suspend Resume Stop Restart Delete

Last Refreshed at 11:45:26 AM

For more information on the crawler monitor, see the KM administration guide.

### Using the TREX Monitor

As soon as the crawling process has finished, you can check the current status of indexing using the TREX monitor.

### Displaying Queues for an Index

In the left-hand side of the window, choose *Display Queues*. Now select the index required in the *Index ID* field. You see an overview of the current processing status of all documents that are in the index queue.

You see, for example, the number of documents that were already prepared for indexing.

Status	Index ID	Queue Status	Timestamp	Delayed	To Be Preprocessed	Preprocessing	Prep. Failed	To Be Transmitted	Transmitting
	how-to-index	Idle	2003-02-04 12:27:30	0	0	0	12	48	0

### Displaying Entries in Queues

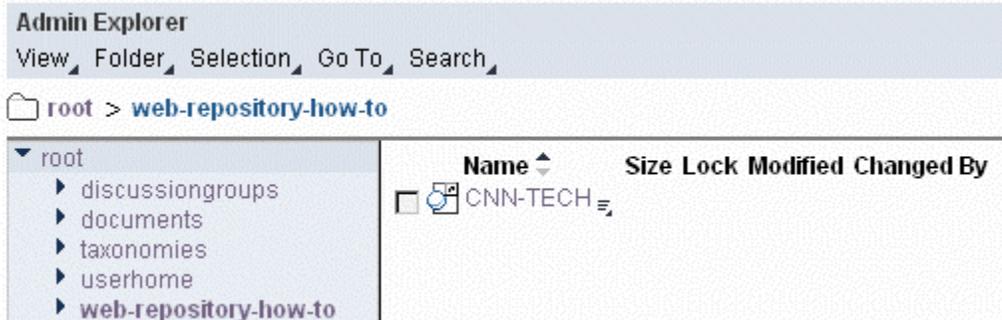
To see an overview of the files contained in an index queue, choose *Display Queue Entries* from the left-hand side of the screen.

The status *OK* means that the indexing of a document has been completed.

Doc ID	Return Code	Return Message	Document Status	Action	Retry Count
/web-repository-how-to/CNN-TECH/2003/TECH/ptech/01/30/dui.science.ap/index.html	0	No error	To Be Transmitted	Index	0
/web-repository-how-to/CNN-TECH/TECH/interne/Archive	0	No error	To Be Transmitted	Index	0
/web-repository-how-to/CNN-TECH/2002/TECH/1/26/thin.wired.id.the/index.html	0	No error	To Be Transmitted	Index	0
/web-repository-how-to/CNN-TECHWEATHER	0	No error	To Be Transmitted	Index	0

### 3 Result

Open the explorer of the Knowledge Management platform. You are now in the *web-repository-how-to* repository that contains the Web site *CNN-TECH*.



When you click on the Web site, it opens in a new window. You can now access the linked documents on the Web site.



Links that are invoked using Javascript cannot be opened.



When you carry out a search, documents from the CNN Web site are also listed in the results list.

[Show Options](#)

**Search Results For** *web attack*

**1-10** [More](#)

 [CNN.com - Iraq war sparks tit-for-tat hacker attacks - Mar. 29, 2003](#) 

83% 4/28/03 9:33:15 AM

Pro-and-anti Iraq war protesters have been making their point by hacking into Web sites in a display of cyber activism, rather than with the traditional can of spray paint or placard.

 [CNN.com - eBay's PayPal accused of violating Patriot Act - Apr. 1, 2003](#) 

63% 4/28/03 9:33:16 AM

A federal prosecutor has alleged eBay Inc. unit PayPal violated a 2001 anti-terror law aimed at fighting money laundering when it provided payment services to online gambling companies, the Web auctio...

 [SAP - SAP to Provide Free IT Consulting Services to Businesses Affected by Recent Terrorist Attacks and Commits \\$3 Million \(U.S.\) For Aid for Victims' Families](#) 

## Copyright

© Copyright 2003 SAP AG. All rights reserved.

No part of this publication may be reproduced or transmitted in any form or for any purpose without the express permission of SAP AG. The information contained herein may be changed without prior notice.

Some software products marketed by SAP AG and its distributors contain proprietary software components of other software vendors.

Microsoft®, WINDOWS®, NT®, EXCEL®, Word®, PowerPoint® and SQL Server® are registered trademarks of Microsoft Corporation.

IBM®, DB2®, DB2 Universal Database, OS/2®, Parallel Sysplex®, MVS/ESA, AIX®, S/390®, AS/400®, OS/390®, OS/400®, iSeries, pSeries, xSeries, zSeries, z/OS, AFP, Intelligent Miner, WebSphere®, Netfinity®, Tivoli®, Informix and Informix® Dynamic Server™ are trademarks of IBM Corporation in USA and/or other countries.

ORACLE® is a registered trademark of ORACLE Corporation.

UNIX®, X/Open®, OSF/1®, and Motif® are registered trademarks of the Open Group.

Citrix®, the Citrix logo, ICA®, Program Neighborhood®, MetaFrame®, WinFrame®, VideoFrame®, MultiWin® and other Citrix product names referenced herein are trademarks of Citrix Systems, Inc. HTML, DHTML, XML, XHTML are trademarks or registered trademarks of W3C®, World Wide Web Consortium, Massachusetts Institute of Technology.

JAVA® is a registered trademark of Sun Microsystems, Inc.

JAVASCRIPT® is a registered trademark of Sun Microsystems, Inc., used under license for technology invented and implemented by Netscape.

MarketSet and Enterprise Buyer are jointly owned trademarks of SAP AG and Commerce One.

SAP, SAP Logo, R/2, R/3, mySAP, mySAP.com and other SAP products and services mentioned herein as well as their respective logos are trademarks or registered trademarks of SAP AG in Germany and in several other countries all over the world. All other product and service names mentioned are trademarks of their respective companies.