

# Documentum DQL How To Configure

## On this page

- [Step 1. Launch Aspire and open the Content Source Management page](#)
- [Step 2. Add a new content source](#)
  - [Step 2a. Specify basic information](#)
  - [Step 2b. Specify the connector information](#)
  - [Step 2c. Specify workflow information](#)
- [Step 3: Initiate a full crawl](#)
  - [During the Crawl](#)
- [Step 4: Initiate an incremental crawl](#)
- [Group Expansion](#)

? Unknown Attachment

## Step 1. Launch Aspire and open the Content Source Management page

Launch Aspire (if it's not already running).

See:

1. [Launch Control](#).
2. Browse to: <http://localhost:50505>. For details on using the Aspire Content Source Management page, see [Admin UI](#).

## Step 2. Add a new content source

To specify exactly which shared folder to crawl, we will need to create a new "Content Source".

? Unknown Attachment

To create a new content source:

1. From the Content Source, select **Add Source**.
2. Select **Connector**.

### Step 2a. Specify basic information

? Unknown Attachment

In the **General** tab in the **Content Source Configuration** window, specify basic information for the content source:

1. Enter a content source name in the **Name** field.
  - a. This is any useful name that you decide is a good name for the source. It will be displayed in the content source page, in error messages, etc.
2. Click on the **Scheduled** list and select one of the following: *Manually, Periodically, Daily, Weekly or Advanced*.
  - a. Aspire can automatically schedule content sources to be crawled on a set schedule, such as once a day, several times a week, or periodically (every N minutes or hours). For the this tutorial, you may want to select **Manually** and then set up a regular crawling schedule later.
3. Click on the **Action** list to select one of the following: *Start, Stop, Pause, or Resume*.
  - a. This is the action that will be performed for that specific schedule.
4. Click on the **Crawl** list and select one of the following: *Incremental, Full, Real Time, or Cache Groups*.
  - a. This will be the type of crawl to execute for that specific schedule.

After selecting Scheduled, specify the details, if applicable:

- *Manually*: No additional options.
- *Periodically*: Specify the **Run every** option by entering the number of "hours" and "minutes".
- *Daily*: Specify the **Start time** by clicking on the hours and minutes lists and selecting options.
- *Weekly*: Specify the **Start time** by clicking on the hours and minutes lists and selecting options, and then selecting the check boxes to specify days of the week to run the crawl.
- *Advanced*: Enter a custom CRON Expression (e.g. 0 0 0 ? \* \*)



You can add more schedules by selecting the **Add New** option, and rearranging the order of the schedules.



If you want to disable the content source, clear the **Enable** checkbox. This is useful if the folder will be under maintenance and no crawls are desired during that period of time.



Real Time and Cache Groups crawl will be available depending on the connector.

## Step 2b. Specify the connector information

In the **Connector** tab, specify the connection information to crawl Documentum.

1. Enter the



m docbase url you want to crawl (Format: `dcfm://<docbroker-server>:<docbroker-port>/<docbase>`).

2. Enter the username (*aspire\_crawl\_account*).
3. Enter the user's password.
4. Enter the location of the *dfc.properties* file. Make sure the *dfc.properties* file correctly points to the *dfc.keystore* in the property: *dfc.security.keystore.file*.
5. Check if you want Error Tolerant: Checked the option if you want to index only metadata and ignore issues during extract content phase.
6. Use RenditionType option. if selected, you have to provide a list of renditions that you want to index. During fetching the document content the first matching rendition from the list will be used provided the document has this type of rendition. If the document has no rendition type other than the default or doesn't match with any of the specified renditions, the connector will use the default.
7. Metadata attributes option. When checked you have to provide the list of the metadata attributes you want to index. If not selected all document attributes will be used. All those attributes appear in the connector specific part of the indexed document.
8. Enter the webtop URL. This URL will be suffixed with the object ID and used as a value of the *displayUrl* element in the resulting metadata. For example: `http://server-name:port/webtop/component/dri?objectId=`
9. Enter a DQL SELECT statement for full crawl. This statement is supposed to query the table with documents - i.e. **dm\_document** and provide the set of IDs for further processing.  
You must use exactly these two mandatory fields in SELECT - **r\_object\_id** and **i\_chronicle\_id**.  
There is also the parameter **{SLICES}** to be used as a part of WHERE clause. If used, the full crawl select would be internally run as 16 other selects (e.g. `<fullselect> AND r_object_id LIKE '%0' - <fullselect> AND r_object_id LIKE '%f'`).  
The purpose of this is to provide an option for parallel processing in the scan phase.
10. Enter DQL SELECT statements for the incremental crawl. In those statements, the parameter **\$(crawlTimeStamp)** should be used. This parameter will be expanded by Aspire at the start of crawling. The time of the last crawl will be used to track down the changes - e.g. `$(crawlTimeStamp)`  
-> `date('09/22/2016 12:00:00','mm/dd/yyyy hh:mi:ss')`:
  - a. The incremental crawl DQL SELECT statement is for picking up "adds" and "updates" with the help of the *r\_modify\_time* attribute in the main table. The *r\_object\_id* and *i\_chronicle\_id* fields are mandatory. This query should be exactly the same as the query for the full crawl plus the parameter `$(crawlTimeStamp)`
  - b. The Audit ACL updates the DQL SELECT statement for tracking ACL changes. The table for querying is **dm\_audittrail**. The fields *r\_object\_id* and *chronicle\_id* are mandatory. The parameter `$(crawlTimeStamp)` should be used.
  - c. The Audit deletes the DQL SELECT statement for tracking deletes (documents or versions of documents). The table for querying is **dm\_audittrail**. The fields *r\_object\_id* and *chronicle\_id* are mandatory. The parameter `$(crawlTimeStamp)` should be used.
  - d. The Audit safeguard SELECT statement is for checking that the IDs of documents (retrieved by the use of a previous audit DQL selects) exists in document table. We can check that the chronicle ID really exists in the document table and that it also exists in the requested WHERE set used in the incremental DQL. Otherwise, we can track all deletes based on the date. However, we would not be sure if the particular audit event is relevant. The *i\_chronicle\_id* field is used and the WHERE part is similar to that in the incremental DQL.  
The parameter **\$(auditChronicleId)** will be expanded and the *chronicle\_id* value from the audit table will be used. The field *r\_object\_id* is mandatory, since this query also serves as a translator between *chronicle\_id* (used in audit table) and *r\_object\_id* (from the document table).
  - e. Select the **Delete audit items** check box if you want Aspire to delete the already processed row in the audit table.
11. Enter the max file size. Any file larger than this size will be ignored by the connector. *Unlimited* includes all files.

12. Select other options as needed
  - a. **Include/Exclude patterns:** This should be handled in WHERE clause of DQL statements.
  - b. **Non-text document filtering:** You should only use the Regex file option for identifying non-text files. The Documentum attribute **a\_content\_type** is used as a non-text filter field.

## Step 2c. Specify workflow information

? Unknown Attachment

In the **Workflow** tab, specify the workflow steps for the jobs that come out of the crawl. Drag and drop rules to determine which steps an item should follow after being crawled. You can use these rules to specify where to publish the document, or which transformations on the data are needed before sending it to a search engine. See [Workflow](#) for more information.

1. For this tutorial, drag and drop the *Publish To File* rule found under the *Publishers* tab to the **onPublish** Workflow tree.
  - a. Specify a **Name** and **Description** for the Publisher.
  - b. Click **Add**.
2. Select **Save** and **Done** to return to the Home Page.

## Step 3: Initiate a full crawl

Now that the content source is set up, the crawl can be initiated.

1. Select the **Full** crawl type option. (The default is Incremental.) The first time it will work like a full crawl. After the first crawl, select **Incremental** to crawl for any changes done in the repository.
2. Select **Start**.

### During the Crawl


During the crawl, you can do the following:

1. Select **Refresh** on the **Content Sources** page to view the latest status of the crawl.  
The status will show **RUNNING** while the crawl is going, and **CRAWLED** when it is finished.
2. Select **Complete** to view the number of documents crawled so far, the number of documents submitted, and the number of documents with errors.  
If there are errors, you will get a clickable Error that will take you to a detailed error message page.

## Step 4: Initiate an incremental crawl

1. If you want to process only content updates from the Documentum (documents that are added, modified, or removed), then select **Incremental** instead of **Full**. The connector will automatically identify only changes which have occurred since the last crawl.

If this is the first time that the connector has crawled, the action of the Incremental option depends on the exact method of *change* discovery. It may perform the same action as a Full crawl and crawl everything, or it may not crawl anything. Thereafter, the Incremental button will only crawl updates.

 Statistics are reset for every crawl.

## Group Expansion

Group expansion configuration is done on the "Advanced Connector Properties" of the Connector tab.

1. Select the **Advanced Configuration** check box to enable the advanced properties section.
2. Scroll to **Group Expansion** and select the check box.
3. Add a new source for each repository from which you want to expand groups from. (You'll need administrator rights on all of them to be able to do this.)
4. Set the default domain, user name, and password of the crawl account.
5. Set a schedule for group expansion refresh and cleanup.
6. As an optional setting select the **Use external Group Expansion** check box to select an LDAP Cache component for LDAP group expansion. See more info on the LDAP Cache component at [LDAP Cache](#).