

User Manual

Printer friendly version

Contents

1. Introduction
2. Selective Harvests
 1. Creating/editing a selective harvest
 1. Easy creation of non existing domains
 2. Event harvest
 3. Adding seeds to an event harvest
3. Snapshot Harvests
 1. Creating/editing a snapshot harvest
4. Domains
 1. Creating Domains
 2. Finding Domains
 3. Editing Domains
 1. Editing configurations
 2. Editing seed lists
 3. Editing crawlertraps
 4. Domain statistics
 5. Alias summary
5. Schedules
6. Global Crawler Traps
7. Heritrix GUI Access
8. Harvest History
 1. All jobs
 2. History of a harvestdefinition
 3. History of a domain
 4. Details on a job
 5. Details on a terminated job
9. Harvester Templates
 1. Download
 2. Upload
10. Quality Assurance
11. System State
12. Bit Preservation
 1. Missing Files
 2. Checksum Errors
13. Alternative Ways to Get Data Out

1. Introduction

The NetarchiveSuite is a software toolset built to be operated by librarians with a

minimum of training in basic internet concepts, webcrawling methods and usage of this system. It's built to handle the national obligations of preserving the internet in Denmark and its design and architecture is much influenced by the Legal Deposit Law stating the three different harvest types:

- snap shot harvests of the national domain
- selective harvests of selected websites
- event harvests in connection with national events

It's also designed to be distributed across the world due to the fact that the danish installation is driven by the two national libraries: The State and University Library, Aarhus and the Royal Library in Copenhagen.

The basic concept in the NetarchiveSuite harvesting module is the notion of domains.

A domain has a two part name [host].[top-level-domain] (e.g. netarchive.dk) or is an IP-number. What is considered a top-level domain is configurable. For most countries it makes sense that the top-level domain is simply the country code (like .dk or .fr), while for others it makes sense to go one level further down (like .co.uk).

A domain can hold multiple so called configurations. A configuration describes how to harvest the domain or a part of the domain. So one configuration could harvest the whole domain (used by the snapshot functionality) and other configurations could take different minor parts of the same domain (for the selective / event harvest).

A configuration consists basically of two things

- A harvester template (predefined templates for the Heritrix web crawler)
- A number of seedlists to use with that template

A domain will always have a default configuration (selectable) and that configuration will be used when starting a snap shot harvest. The snap shot harvest therefore takes all domains known in the database.

edit

2. Selective Harvests

selective_harvests.png

The front page by default shows the list of selective harvests.

You can [**Activate**] an inactive harvest definition and [**Deactivate**] an active harvest definition. If you deactivate a running harvest, the system will finish the running jobs.

Click on [**Edit**] to change an existing harvest definition or [**Create new harvest definition**]

Click on [**History**] if you wish to trace back all the jobs from former finished harvests.

2.1. Creating/editing a selective harvest

selective_harvest_edit.png

Create a new selective harvest definition by pressing [**Create new selective harvestdefinition**] from the frontpage.

Give the harvestdefinition a recognizable harvest name – you can not change it later. If necessary add a comment.

Choose a schedule from the dropdown list.

Now you can add domains to the harvestdefinition.

Write the name of the domains you want to add in the box “**Enter domain(s) to add to the harvest here**” and click on [**Add domains**].

The added domains will appear in the column “Domain”.

For each added domain, choose the wanted configuration from the dropdown list for each domain. Press [**Save**] to save the harvestdefinition.

The scheduling of selective harvest definitions can be overridden by filling out the input field **Override with new date**. Simply set the date to whenever you wish the harvest definition to run next time. The scheduling of the harvest definition will continue from that point in time.

2.1.1. Easy creation of non existing domains

selective_harvest_non_existing_domain.png

When adding a domain that is not existing in the database you are warned with *The following domains are unknown and were not added*. You can simply add the unknown domains to the database and your harvestdefinition by clicking [**Create and add to the harvestdefinition**]

2.2. Event harvest

Event harvests are treated almost the same as selective harvests in the system. The only difference is a power-adding of domains function. This could be used for selective harvests as well but was developed for event harvesting definitions where the operator must fill in larger number of URLs without having to edit configurations and seedlists on all those domains.

2.3. Adding seeds to an event harvest

event_harvest_seeds.png

Use [*Add seeds*]. Enter identified start-URLs covering the event in the “Enter seeds:” box. In “Max number of bytes per domain” enter preferred max number, e.g. 1000000000. Select a harvest template with the “Harvest template” drop down box.

All seeds will use the same template, so to harvest different seeds with different templates you need to add them bunch by bunch for each template you need for your event harvest.

Pressing [*Insert*] starts the power-adding function. This function runs through the entered seeds one by one and does the following with each seed:

1. Finds the domain from which the seed derives
2. Creates a seedlist with the name of the harvestdefinition and the template as seedlist-name
3. Creates a configuration with the name of the harvestdefinition and the template as configuration-name. And select the seedlist from (2) to use with the new configuration. If the seedlist to create in (2) or the configuration to create in (3) already exist (If the power-adding function has been used before with other seeds from the same domain in the same event harvest) the system will only add the new URLs to the existing seedlist.

You can also use [*Add seeds from a file*]. This allows you upload a file with the seeds instead of entering the seeds in a text field. Otherwise the functionality is the same.

event_harvest_seeds_from_file.png

edit

3. Snapshot Harvests

snapshot_harvests.png

On [***Snapshot Harvests***] new snapshot harvests are started, harvesting all domains known to the system in their default configurations. An overview of all snapshot harvests is also provided.

[***Create new harvestdefinition***] opens the template below.

3.1. Creating/editing a snapshot harvest

snapshot_harvest_edit.png

This page is used to define name and size (max. bytes per domain) of the harvest. It is now possible to use number of objects as harvest limits, as well as the size in bytes. The default object limit for harvests if using object limits rather than bytelimits. -1 means unlimited.

It is recommended to systematize the naming for clarity, e.g. 2007-1, 2007-2 etc.

The size of the harvest can be defined in two ways: at the harvest definition [***Snapshot Harvests***] or at the configuration of the single domain. It will always be the lower size limit stopping the harvesting of a domain.

Comments can freely be added.

Snapshot harvests can be based on previous snapshots in the sense that it can be limited to only harvest domains that hit the max number of bytes limit in a previous harvest.

The domains completely finished (not hitting the max number of bytes limit – either on the configuration level or on the snapshot harvest level) in the first harvest will not be included in the second. Domains included in harvests which were aborted through the Heritrix GUI or otherwise stopped uncleanly (for example by a crash of a harvester machine) will also not be included.

All other domains will be harvested from the beginning in the second harvest.

[***Save***] saves the harvest definition and returns to [***Snapshot harvests***].

After defining a snapshot harvest the harvest is activated with the [***Activate***] button on the snapshot frontpage. Harvest will not start until you press [***Activate***]. Status then changes to 'Active'.

[***Deactivate***] is not relevant in Snapshot Harvests because they only run once. By [***Edit***]

the Snapshot Definition can be changed but only before activation. Parameters changeable are size, commentary and if previous harvest startpoint should be used. The name can not be changed.

[*History*] provides an overview of the specific harvest: see User Manual 3.12/Harvest History

edit

4. Domains


4.1. Creating Domains

create_domains.png

The [*Create domain*] is used for creating new domains in the system. It is possible to create a single domain as well as list of domains. It is also possible to import domains from a file.

To create single domains enter domain names in the text box and press [*Create*].

To bulk create domains from a file select the file from your local computer with [*Browse*] and press [*Ingest*]. The file must be a simple list of domain names – one at each line. The file must be UTF-8 encoded if it contains special characters.

New domains get a default configuration when created (with the defaultorderxml template and a default maximum number of bytes). New domains also get a defaultseedlist when created. The frontpage is defined as `www.domainname` – e.g. '  `www.netarchive.dk`' for the domain 'netarchive.dk'

Already existing domains in the system will not be recreated.

4.2. Finding Domains

find_domains.png

[*Find Domain(s)*] is used to find domains existing in the system.

Write a domain name in the box (e.g. kb.dk). Searching is done on the complete text string. Press [*Search*].

Left and/or right wildcards with *.

If there are several hits, a list is given of the found domains. If only one hit it leads directly to the Domain page.

If the search for a specific domain results in no hits, you are prompted with the ability to create the domain in the system and by accepting [*Yes*] it leads directly to Domain page for the newly created domain.

4.3. Editing Domains

domain_edit.png

'*Edit domain*' is an overview of a single domain where it is possible to edit the domain's definition in the harvest system.

Free commentary text box.

'Alias of': Here it can be stated if the domain is an alias of another domain – they are identical in content and only one of them should be harvested. Domains marked as an alias will not be harvested within the snapshot harvests. Alias is defined one year at a time and then has to be renewed.

'Configurations': [*New configuration*] and [*Edit*] open a new page: '*Enter/edit configuration*' (see below) 'Seed lists': [*New seed list*] and [*Edit*] opens a new page: '*Enter/edit seed list*'. (see below)

'Crawler traps': [*Show crawler traps*] opens a new text box: '*Crawler traps*' (see below)

[*Show historical harvest information for ...*] opens a new page "**Harvest history for domain....**" (see User Manual 3.12/Harvest History) ['

4.3.1. Editing configurations

configuration_edit.png

'*Enter/edit configuration*' is used to define a new configuration and edit an existing one. A configuration contains information about which Harvest template and Seed lists are used (more than one Seed list can be used - hold down CTRL).

At the creation of a new configuration a name is given that thereafter can not be changed.

Furthermore it is possible to choose between different Harvester templates and maximum number of bytes to be harvested in each harvest of the configuration. At creation the default number of bytes is chosen for each domain. And a default maximum number of objects is set, but can be overwritten.

4.3.2. Editing seed lists

seedlist_edit.png

'Enter/edit seed list' is used to define a new Seed list or to edit an existing one.

At the creation of a new Seed list a name is given that thereafter can not be changed.

In the 'Seeds' text box a list of seeds to be harvested is given. Seeds can be omitted by writing a # prefix, e.g. # <http://www.kb.dk>. This can also be used for comments inside the seedlist – e.g. '#this seed is important'

4.3.3. Editing crawlertraps

crawlertrap.png

A crawlertrap is a path followed blindly by the harvester which in principle can continue forever. This typically could be a calendar.

To avoid crawlertraps on a domain, the administrator can state parts of URLs that should never be harvested (in any configuration). Matching URLs are omitted in all harvests of the domain and in other domains harvested in the same job. So be careful not to give too general statements that could potentially omit things on other domains (perhaps always include the domainname itself in the statement).

The string of text must be stated as a 'regular expression'.

4.4. Domain statistics

domain_statistics.png

The domain statistics page will give you information about number of subdomains for each unique Top level domain known in the system. IP-numbers will be counted separately.

The number in the "Number of subdomains" column is clickable and will do a search for all domains matching that Top level domain. This is only applicable to Top level domains with a limited number of subdomains since the matching domains will

be listed on one page – and that page will get very long if the system contains hundreds of thousands of domains.

4.5. Alias summary

alias_summary.png

The alias summary page gives an overview of the domains marked as aliases of other domains in the system. Both domain names are clickable and will open the domain page for the clicked domain.

The “Expires” column shows when the alias expires (12 month after they are marked). The mark does not disappear after 12 month in the database but the “Overview of Aliases” page will show the “expired” ones in the top.

To renew an alias for another 12 month one is currently forced to open the domain page of the marked domain (the “Domain” column) – select “renew alias” and press “Save”

edit

5. Schedules

schedules.png

Schedules are only applied on selective and event harvests.

A schedule defines a harvesting frequency. The minimum entity is one hour. It is possible to choose an automatically fixed start and/or end time for a specific harvest.

It is possible to create an infinite number of schedules.

For a new schedule click on [Create new schedule]. And to edit an existing schedule press [Edit].

schedules_edit.png

Give the schedule an easily recognizable name – note that it can't be changed once saved. If necessary, add a comment.

Fill in the frequency and – if necessary - time of the day for the harvest to run. In the drop down menu you have the choice between hours, days, weeks and months.

Changing “days” switches the “Time of day” so that:

- *“hours” lets you choose a specific minute of the hour*
- *“days” lets you choose a specific time of day*
- *“weeks” lets you choose a specific day of the week*
- *“months” lets you choose a specific day of the month*

After selecting the frequency you must select “Start at the earliest” which could either be as soon as possible (default) or at a specific date and time.

The last thing to determine is how long this schedule should go on. The default for the duration of a schedule is forever. It is also possible to choose an end date and a certain number of harvests to perform. This allows you to define schedules that will only run in a shorter period – e.g. in connection with an event harvest where the date range in which to harvest is predefined.

edit

6. Global Crawler Traps

A crawler trap is any sequence of webpages which a crawler can blindly and endlessly follow without harvesting any new information. A common example is a calendar system with hyperlinks to subsequent or previous dates. Crawler traps can be avoided by specifying (as regular expressions) URLs which the crawler is to ignore. In NetarchiveSuite one can specify crawler traps either per-domain or globally. This section describes the management of global crawler traps.

A list of crawler traps is just a plain text file containing crawler-trap regular expressions one-per-line. Lists may be active or inactive. When NetarchiveSuite creates a new job for any harvest, all crawler traps for all active lists (excluding duplicates) are added to the crawl template for that job.

GlobalCrawlerTraps_1.png

To upload a list of global traps, first click on the [Edit] link and fill in a name and description for the list of crawler traps and the path where the file containing the crawler trap expressions is to be found. You can also choose whether the list should be initially active or inactive. Click [Create] to upload the list.

GlobalCrawlerTraps_2.png

A list may be made active or inactive by clicking on the [Activate] and [Deactivate] buttons. Lists may also be viewed (via the [Retrieve] button), deleted, or edited. Note that the retrieved version of a crawler trap list may differ from the original uploaded version because any duplicates in the original are removed during upload and the order of the lines in the retrieved version will not be the same as in the original file. The [Edit] actions allow for uploading of a new version of the list.

GlobalCrawlerTraps_3.png

A side effect of using global crawler trap lists is that the database will grow more rapidly as the modified crawl template, including all the active crawler traps, is stored for every job.

edit

7. Heritrix GUI Access

It is possible - only while a job is running - to access the Heritrix user interface on the harvester machine. Start a browser on the harvester machine and use the port specified, e.g. <http://my.harvester.machine:8090>. The port is defined by the setting `settings.harvester.harvesting.heritrix.port`. Enter the administrator name e.g. "admin" and password e.g. "adminPassword" as set in the `settings.harvester.harvesting.heritrix.adminName` and in `settings.harvester.harvesting.heritrix.adminPassword` settings.

See in the installation manual how you change settings.

In the Heritrix GUI you can e.g. pause, stop or restart a job.

heritrixGUI.gif

edit

8. Harvest History

8.1. All jobs

[Harvest Status] in the left menu by default shows a chronological list of all jobs ever harvested with status Started in ascending order. The same does [All jobs]

all_jobs.png

If information is wanted for jobs with other statuses (or All statuses) or other sort order, then this can be specified in the combo-boxes in the top of the page and then activated by clicking the Show button.

For each job the page shows information about the job and its status as well as information about errors (harvest errors or upload errors) and number of configurations in the job.

Chose [Run number] if you want to check details on a specific run of that harvestdefinition – note that a run can consist of multiple jobs.

Chose [Harvest name] if you want to check details on the history of a specific harvest definition.

Chose [JobID] if you want to check details on a specific job.

In case of Harvest errors, a [Restart] button will appear and the operator can choose to resubmit that specific job to be harvested again.

failed_desc.png

When resubmitting a failed job, the status will say 'Resubmitted' and refer to the new resubmitted job.

resubmit_job.png

8.2. History of a harvestdefinition

history_harvestdefinition.png

The history page for a harvestdefinition is the same as you can reach from the frontpage with the [History] buttons.

This history page gives you further information for each run of the harvestdefinition: Start time, End time, number of bytes harvested and number of documents harvested.

The page also show how many jobs each run consists of and how many of these that failed and eventually got resubmitted.

8.3. History of a domain

history_domain.png

If you want to see all the jobs connected to a specific domain, click on [All jobs per domain] and search for the domain name.

You will get a chronological list of the harvest definitions including the chosen domain.

This page gives the same history information as the other two history pages and further more gives a “Stopped due to” information. This column will show the operator if a harvest was stopped unexpectedly or if the harvest hit the max-bytes limit for the chosen domain or if the harvest was stopped because of an error on the harvester machine.

8.4. Details on a job

job_details.png

Clicking on a jobID on any of the harvest history pages will give you a very detailed report on the job.

This page gives all the information available about the job itself (e.g. max-bytes limit) and about the single domains included in the job.

Furthermore the page shows the complete seedlist used with the job and the complete “Harvest order template” as well as detailed error information in case of errors. The two latter is mainly for advanced users debugging specific crawls where things didn't go as expected.

8.5. Details on a terminated job

job_details_terminated.png

If you terminates a running job in the Heritrix GUI, you can view Job Details and see that the job is stopped due to "Harvesting aborted".

If it is a bigger (snapshot) harvest, that includes several jobs, all the finished jobs, will appear as "Done".

And only the ones that was actually stopped will appear as stopped due to "Harvesting aborted" for some domains.

edit

9. Harvester Templates

edit_harvest_templates.png

The [Edit Harvest Templates] is used for managing the harvester templates. It enables you to both download and upload templates from/to the system database.

9.1. Download

The download part lets you view existing templates as either plain text or XML in the browser window or download existing templates to your local computer.

Select the template you want to view/download in the first select box, select the method in the second and press [Retrieve]

9.2. Upload

In the Upload section you can either update an existing template or create a new one.

To update an existing template, select the template to update in the first select box and browse for a file on your local computer with the [Browse] button.

By pressing [Replace harvesttemplate with file from your own harddrive] you will overwrite the chosen template in the database.

To create a new template give the new template a name in the “Template Name” box and select a file from your local computer with [Browse].

By pressing [Create a new harvesttemplate using a file from your own harddrive] you will add the new template with the given name to the database.

When using the upload functions the uploaded files will be checked against certain rules to ensure that the templates contain specific elements used by the NetarchiveSuite system.

edit

10. Quality Assurance

viewerproxy_status.png

Quality assurance is done by browsing the archive for selected domains. If something is missing on the pages the system can be set to automatically collect all the missing URL's for later transfer to the harvesting system. Before doing Quality Control you need to setup your browser to use a proxyserver (see Quick Start Manual)

It is suitable to investigate one domain at a time (unless several domains are included in the same website-complex).

[Start collecting URL's] Hereby starts the collection of URL's. The Current Viewerproxy status textbox shows if the system collects URL's or not – and how many URL's are currently collected.

[Stop collecting URL's] Collection of URL's can be stopped at any time.

[Clear collected URL's] The list of URL's can be cleared at any time e.g. when investigating a new domain starts. NB! This function can not be undone.

[Show collected URL's] The list of URL's can be viewed at any time. The list can be copied and manually be added to relevant Seed lists for the relevant domains in the harvesting system.

edit

11. System State

The [Systemstate] pages lets the operator monitor the entire system (all machines and applications) from one central point.

systemstate_all.png

The initial view is the last log-message from every machine and every application.

This can be narrowed down to single “Machine” / “Application” / “Instance id” / “Priority” / “Use Replica” and extended to not only show the last log-message but the last 100 log-messages (don't do that for the initial view of everything).

To narrow the view press either of the links on the page – when narrowed the view can be extended again with the [Show all] buttons that will dynamically appear in the headlines of the table.

It is from 3.10 also possible to remove an application from the system. Be careful about this feature, because removing a running application, will make it disappear. The new column and button is added to the right: [Remove Application]

The following is a view of one harvester instance on one specific machine (narrowed down by application):

systemstate_one_app.png

If you load this page just while a harvester instance is restarting you might get a JMX-error. The same thing will happen if one of the configured applications does not run or does not respond. So the system state will in some sense also discover non functional applications.

edit

12. Bit Preservation

bitpreservation.png

The [Bitpreservation] interface lets you control active checks of the status of the underlying bitarchive. This only applies if your installation uses the NetarchiveSuite bitarchive application.

The interfaces lets you initiate two types of checks on every copy of the files in the archive: Filestatus and Checksum status.

In the example on the screen dump there are more bitarchive instances – e.g. SBN and KBN.

The [Update] buttons let you update the status for both files and checksums for both bitarchive instances. The page will give you the Filestatus as

- *Number of files*
- *Missing files*
- *Last updated at*

and the Checksum status as

- *Number of files with error*
- *Last updated at*

The Filestatus checks are rather fast because only the existence of the files are

checked whereas the Checksum status checks can take days/weeks for larger archives depending on the number of CPUs and the IO-speed of your hard drives.

12.1. Missing Files

missing_files.png

If files are missing on one instance of the bitarchive a [Show missing files] button will appear right next to the line with the number of missing files

For each missing file you can select “Get info”. With the “Change the infobox for”-field in the bottom of the screen you can select a number of files in one operation.

Pressing [Execute] makes the system get a fresh status on the files and their checksums from both copies (out of which one is missing) and from the administrative system so that there are always three checksums available for each file.

missing_files_info.png

If the two “remaining” checksums (On the screen dump 32 byte long) are identical the system allows you to add the missing file to the bitarchive instance that had lost it.

This requires you to click “Add to archive” and then press [Execute]. Marking files for addition can be done for a number of files in one operation.

12.2. Checksum Errors

checksum_error.png

In case of a checksum error this error can be corrected through the interface. To replace a bad file you need to type a security password and press [Replace the file in bitarchive replica XX]. The bad file will not be completely removed but moved to an 'attic' directory on the bitarchive server holding the bad file.

edit

13. Alternative Ways to Get Data Out

There are alternative ways to get data out of the bitarchive, e.g. by run of batch programs on the bitarchive replicas. However, there are no explicit user interface to

these tools and some technical skills are required to use them. This is the reason why these tools are described in the Additional Tools Manual under the tools in the Archive Module.

edit

User Manual 3.12 (last edited 2010-08-16 10:24:50 by localhost)