

# User Manual

Printer friendly version

Contents	
1.	Introduction
2.	Selective Harvests
1.	Creating/editing a selective harvest
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.	Heritrix GUI Access
7.	Harvest History
1.	All jobs
2.	History of a harvestdefinition
3.	History of a domain
4.	Details on a job
8.	Harvester Templates
1.	Download
2.	Upload
9.	Quality Assurance
10.	System State
11.	Bit Preservation
1.	Missing Files
2.	Checksum Errors

## 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

## Selective Harvests

Dansk English

**Menu**

- Definitions
- Selective Harvests
- Snapshot Harvests
- Schedules
- Find Domain(s)
- Create Domain
- Domain Statistics
- Alias Summary
- Edit Harvest Templates
- Harvest status
- Bitpreservation
- Quality Assurance
- Systemstate

## Snapshot Harvests

Harvest definition	Max bytes	Status	Commands		
my_first_harvest	100000	Active	Deactivate	Edit	History
my_second_harvest	100000	Inactive	Activate	Edit	History

Create new snapshot harvest definition

NetarchiveSuite Version: 3.6.0 status RELEASE, TEST2

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 a new harvest definition**]

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

## Creating/editing a selective harvest

Dansk English

**Menu**

- Definitions
- Selective Harvests
- Snapshot Harvests
- Schedules
- Find Domain(s)
- Create Domain
- Domain Statistics
- Alias Summary
- Edit Harvest Templates
- Harvest Status
- Bitpreservation
- Quality Assurance
- Systemstate

## Selective Harvest

Harvest name:

Comments:

Schedule:

The harvestdefinition my\_fifth\_harvest is inactive. If activated, it will run again on Jun 13, 2008 9:06:13 AM  
Override with new date  (format: DD/MM YYYY hh:mm)

Domain	Choose configuration	Remove from harvest
netarkivet.dk	defaultconfig	Remove
netarchive.dk	defaultconfig	Remove

Enter domain(s) to add to the harvest here:

Event harvest:

NetarchiveSuite Version: 3.6.0 status RELEASE, TEST2

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.

## 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.

## Adding seeds to an event harvest

The screenshot shows the 'Event harvest: my\_fifth\_harvest' configuration page. On the left is a 'Menu' sidebar with categories: Definitions (Selective Harvests, Snapshot Harvests, Schedules, Find Domain(s), Create Domain, Domain Statistics, Alias Summary, Edit Harvest Templates), Harvest status, Bitpreservation, Quality Assurance, and Systemstate. The main content area has a header 'Event harvest: my\_fifth\_harvest' and a sub-header 'Enter seeds:'. Below this is a large empty text box for entering seeds. At the bottom of the main area, there are two input fields: 'Max number of bytes per domain:' with the value '1000000000' and 'Harvest template:' with a dropdown menu showing '3levels\_orderxml'. An 'Insert' button is located below these fields. The bottom of the page features a blue footer bar with the text 'NetarchiveSuite Version: 3.6.0 status RELEASE, TEST2'. The top right corner of the page shows 'Dansk English'.

Use **[Add seeds]**. Enter identified start-URLs covering the event in the “Enter seeds:” box. 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.

edit

## Snapshot Harvests

Dansk English

**Menu**

- Definitions
- Selective Harvests
- Snapshot Harvests
- Schedules
- Find Domain(s)
- Create Domain
- Domain Statistics
- Alias Summary
- Edit Harvest Templates
- Harvest status
- Bitpreservation
- Quality Assurance
- Systemstate

## Snapshot Harvests

Harvest definition	Max bytes	Status	Commands		
my_first_snapshot	10000000	Inactive	Activate	Edit	History

Create new snapshot harvest definition

NetarchiveSuite Version: 3.6.0 status RELEASE, PRO

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.

## Creating/editing a snapshot harvest

Dansk English

**Menu**

- Definitions
- Selective Harvests
- Snapshot Harvests
- Schedules
- Find Domain(s)
- Create Domain
- Domain Statistics
- Alias Summary
- Edit Harvest Templates
- Harvest status
- Bitpreservation
- Quality Assurance
- Systemstate

## Snapshot Harvest

Harvest name:

Max number of bytes per domain:

Comments:

**Harvest only domains that were not completely harvest in a previous harvest:**

None	<input type="checkbox"/>
my_first_harvest	<input type="checkbox"/>

NetarchiveSuite Version: 3.6.0 status RELEASE, TEST2

This page is used to define name and size (max. bytes per domain) of the harvest. 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.

All other domains will be harvested from the beginning up till the newly specified limit for the new snapshot.

[*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/Harvest History

edit

## Domains

## Creating Domains

Menu

Definitions

- Selective Harvests
- Snapshot Harvests
- Schedules
- Find Domain(s)
- Create Domain
- Domain Statistics
- Alias Summary
- Edit Harvest Templates

Harvest status

Bitpreservation

Quality Assurance

Systemstate

Dansk English

### Create Domain

Enter the domain or list of domains to be created:

Create

Or import domains from a file on your harddrive

 Browse... Ingest

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 maximum number of bytes at 500Mbytes). New domains also get a defaultseedlist when created. This seedlist contains a link to the frontpage of the domain. 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.

## Finding Domains

Menu

Definitions

- Selective Harvests
- Snapshot Harvests
- Schedules
- Find Domain(s)
- Create Domain
- Domain Statistics
- Alias Summary
- Edit Harvest Templates

Harvest status

Bitpreservation

Quality Assurance

Systemstate

Dansk English

### Find Domain(s)

Enter the name of the domain(s) to find:  Search

(Use of wildcards permitted, e.g., \*.com)

NetarchiveSuite Version: 3.6.0 status RELEASE, TEST2

[**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.

## Editing Domains

'*Edüt 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 100% 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*] Open a new text box: '*Crawler traps*' (see below)

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

### Editing configurations

'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. The default number of bytes is 500Mbytes

### Editing seed lists

The screenshot shows the 'Enter / edit seed list' form in the NetarchiveSuite web interface. The domain is 'netarkivet.dk'. The 'Name' field contains 'defaultseeds'. The 'Seeds' field contains 'www.netarkivet.dk'. There is an empty 'Comments' field. A 'Update seed list' button is located at the bottom of the form. The interface is in Danish, and the version is 3.6.0 status RELEASE, TEST2.

'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'

### Editing crawlertraps

The screenshot shows the 'Crawler traps' form in the NetarchiveSuite web interface. The text box contains the regular expression 'netarchive\.dk\/calendar\.php.\*'. An 'Update Domain' button is located at the bottom of the form. The interface is in Danish, and the version is 3.6.0 status RELEASE, TEST2.


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.

### Domain statistics

Dansk English



**Menu**

- Definitions
- Selective Harvests
- Snapshot Harvests
- Schedules
- Find Domain(s)
- Create Domain
- Domain Statistics
- Alias Summary
- Edit Harvest Templates
- Harvest status
- Bitpreservation
- Quality Assurance
- Systemstate

## Domain Statistics

**Number of registered domains: 19**

Top level domain	Number of subdomains
com	1
dk	18


NetarchiveSuite Version: 3.6.0 status RELEASE, TEST2

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.

## Alias summary

Dansk English



**Menu**

- Definitions
- Selective Harvests
- Snapshot Harvests
- Schedules
- Find Domain(s)
- Create Domain
- Domain Statistics
- Alias Summary
- Edit Harvest Templates
- Harvest status
- Bitpreservation
- Quality Assurance
- Systemstate

## Overview of Aliases

**Existing aliases:**

Domain	Alias of:	Expires
<a href="#">netarkivet.dk</a>	<a href="#">netarchive.dk</a>	Jun 12, 2009 11:11:21 AM
<a href="#">detkongeligebibliotek.dk</a>	<a href="#">kb.dk</a>	Jun 12, 2009 11:13:33 AM

NetarchiveSuite Version: 3.6.0 status RELEASE, TEST2

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

## Schedules

Dansk English



**Menu**

- Definitions
- Selective Harvests
- Snapshot Harvests
- Schedules
- Find Domain(s)
- Create Domain
- Domain Statistics
- Alias Summary
- Edit Harvest Templates
- Harvest status
- Bitpreservation
- Quality Assurance
- Systemstate

## Schedules

**Existing Schedules:**

Once_a_day	<a href="#">Edit</a>
Once_a_month	<a href="#">Edit</a>
Once_a_week	<a href="#">Edit</a>
Once_an_hour	<a href="#">Edit</a>

[Create new schedule](#)

NetarchiveSuite Version: 3.6.0 status RELEASE, TEST2

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*].

Menu

Definitions

- Selective Harvests
- Snapshot Harvests
- Schedules
- Find Domain(s)
- Create Domain
- Domain Statistics
- Alias Summary
- Edit Harvest Templates

Harvest status

Bitpreservation

Quality Assurance

Systemstate

Dansk English

## Edit Schedule

Schedule name:

Comments:

**Perform harvest**

Every  days

Time of day:  Anytime  
 At  :

**Starts at the earliest**

As soon as possible  
 at  (format: DD/MM YYYY hh:mm)

**Continue**

Forever  
 Until  (format: DD/MM YYYY hh:mm)  
 Until  harvests have been done

NetarchiveSuite Version: 3.6.0 status RELEASE, TEST2

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

## 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.



Status as of **May. 2, 2008 12:58:41 GMT** Alerts: no alerts  
 CRAWLING JOBS RUNNING job: 2-2  
 Admin Console 0 jobs pending, 0 completed 3270 URIs in 2m52s (14.2/sec)

[Console](#) [Jobs](#) [Profiles](#) [Logs](#) [Reports](#) [Setup](#) [Help](#)

Crawler Status: **CRAWLING JOBS** | [Hold](#)

<b>Jobs</b>	<b>Memory</b>
Running: 2-2	48377 KB used
0 pending, 0 completed	54208 KB current heap
Alerts: 0 (0 new)	1575808 KB max heap

---

Job Status: **RUNNING** | [Pause](#) | [Checkpoint](#) | [Terminate](#)

<b>Rates</b>	<b>Load</b>
14.2 URIs/sec (19.44 avg)	8 active of 50 threads
294 KB/sec (684 avg)	1 congestion ratio
<b>Time</b>	1969 deepest queue
2m52s elapsed	474 average depth
5m24s remaining (estimated)	
<b>Totals</b>	
downloaded 3270	6164 queued
9442 total downloaded and queued	
115 MB crawled (115 MB novel)	

[Refresh](#)


edit

## Harvest History

### 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]

Dansk English

 <b>Menu</b> <a href="#">Definitions</a> <a href="#">Harvest status</a> All Jobs All Jobs per domain <a href="#">Bitpreservation</a> <a href="#">Quality Assurance</a> <a href="#">Systemstate</a>	Only show job status <input type="text" value="Started"/> in <input type="text" value="Ascending"/> order <input type="button" value="Show"/>													
	<h3>Job Status</h3> <table border="1"> <thead> <tr> <th>Job ID</th> <th>Harvest name</th> <th>Run number</th> <th>Status</th> <th>Harvest errors</th> <th>Upload errors</th> <th>Number of configurations</th> </tr> </thead> <tbody> <tr> <td>1</td> <td>my_first_harvest</td> <td>0</td> <td>Started</td> <td>-</td> <td>-</td> <td>1</td> </tr> </tbody> </table>	Job ID	Harvest name	Run number	Status	Harvest errors	Upload errors	Number of configurations	1	my_first_harvest	0	Started	-	-
Job ID	Harvest name	Run number	Status	Harvest errors	Upload errors	Number of configurations								
1	my_first_harvest	0	Started	-	-	1								

NetarchiveSuite Version: 3.6.0 status RELEASE, TEST

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.


Dansk English

 <b>Menu</b> <a href="#">Definitions</a> <a href="#">Harvest status</a> All Jobs All Jobs per domain <a href="#">Bitpreservation</a> <a href="#">Quality Assurance</a> <a href="#">Systemstate</a>	Only show job status <input type="text" value="All"/> in <input type="text" value="Ascending"/> order <input type="button" value="Show"/>																											
	<h3>Job Status</h3> <table border="1"> <thead> <tr> <th>Job ID</th> <th>Harvest name</th> <th>Run number</th> <th>Status</th> <th>Harvest errors</th> <th>Upload errors</th> <th>Number of configuration</th> </tr> </thead> <tbody> <tr> <td>1</td> <td>my_first_harvest</td> <td>0</td> <td>Done</td> <td>-</td> <td>-</td> <td>16</td> </tr> <tr> <td>2</td> <td>my_fifth_harvest</td> <td>0</td> <td>Failed</td> <td>dk.netarkivet.common.exceptions.IOFailure: Crawl probably interrupted by shutdown of HarvestController <input data-bbox="981 1886 1045 1908" type="button" value="Restart?"/></td> <td>-</td> <td>1</td> </tr> <tr> <td>3</td> <td>my_fifth_harvest</td> <td>0</td> <td>Done</td> <td>-</td> <td>-</td> <td>1</td> </tr> </tbody> </table>	Job ID	Harvest name	Run number	Status	Harvest errors	Upload errors	Number of configuration	1	my_first_harvest	0	Done	-	-	16	2	my_fifth_harvest	0	Failed	dk.netarkivet.common.exceptions.IOFailure: Crawl probably interrupted by shutdown of HarvestController <input data-bbox="981 1886 1045 1908" type="button" value="Restart?"/>	-	1	3	my_fifth_harvest	0	Done	-	-
Job ID	Harvest name	Run number	Status	Harvest errors	Upload errors	Number of configuration																						
1	my_first_harvest	0	Done	-	-	16																						
2	my_fifth_harvest	0	Failed	dk.netarkivet.common.exceptions.IOFailure: Crawl probably interrupted by shutdown of HarvestController <input data-bbox="981 1886 1045 1908" type="button" value="Restart?"/>	-	1																						
3	my_fifth_harvest	0	Done	-	-	1																						

NetarchiveSuite Version: 3.6.0 status RELEASE, TEST

### History of a harvestdefinition

Dansk English



**Menu**

- Definitions
- Harvest status
  - All Jobs
  - All Jobs per domain
- Bitpreservation
- Quality Assurance
- Systemstate

## Harvest history for full harvest my\_first\_harvest

Run number	Start time	End time	Bytes Harvested	Documents Harvested	Total number of jobs	Number of failed jobs	Number of resubmitted jobs
0	Jun 9, 2008 4:08:59 PM	Jun 9, 2008 4:11:37 PM	2,304,814	493	1 <a href="#">Show jobs</a>	0	0

NetarchiveSuite Version: 3.6.0 status RELEASE, TEST2


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.

## History of a domain

Dansk English



**Menu**

- Definitions
- Harvest status
  - All Jobs
  - All Jobs per domain
- Bitpreservation
- Quality Assurance
- Systemstate

## Harvest history for netarkivet.dk

Harvest name	Run number	Job ID	Configuration	Start time	End time	Bytes Harvested	Documents Harvested	Stopped due to
my_fifth_harvest	0	3	defaultconfig	Jun 13, 2008 11:46:16 AM	Jun 13, 2008 11:47:40 AM	11,470,199	124	Domain Completed

NetarchiveSuite Version: 3.6.0 status RELEASE, TEST2


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 pages 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 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.

## Details on a job

Dansk English



**Menu**

- Definitions
- Harvest status
  - All Jobs
  - All Jobs per domain
- Bitpreservation
- Quality Assurance
- Systemstate

## Details for Job 4

Job ID	Type	Harvest name	Run number	Start time	End time	Status	Harvest errors	Upload errors	Object limit	Byte limit
4	Full harvest	Test3	0	Jun 6, 2008 2:28:30 PM	Jun 6, 2008 2:29:35 PM	Done			-1	40,000

### QA job selection

Select this job for QA with viewerproxy

The link above will select this job for the viewerproxy browse index. This will only work if your browser is set up to use the viewerproxy as web proxy.

### Included domains and configurations

Domain	Configuration	Bytes Harvested	Documents Harvested	Stopped due to
raeder.dk	defaultconfig	47,127	13	Domain-config limit reached

### Seed list

www.raeder.dk

### Harvest order template (based on default\_orderxml)

Show harvest template for job 4

NetarchiveSuite Version: 3.6.0 status RELEASE, TEST2

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.

edit

## Harvester Templates

Dansk English

### Menu

- Definitions
- Selective Harvests
- Snapshot Harvests
- Schedules
- Find Domain(s)
- Create Domain
- Domain Statistics
- Alias Summary
- Edit Harvest Templates
- Harvest status
- Bitpreservation
- Quality Assurance
- Systemstate

### Edit Harvest Templates

#### Download

Select one of the following templates:

3levels\_orderxml Show as text Retrieve

#### Upload

Here you can upload a harvesttemplate to replace an existing harvesttemplate:

3levels\_orderxml Gennemse...  
Replace harvesttemplate with file from your own harddrive

Here you can create a new harvesttemplate by uploading it from your harddrive:

Template Name Select file: Gennemse...  
Create a new harvesttemplate using a file from your own harddrive

NetarchiveSuite Version: 3.6.0 status RELEASE, TEST3

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.

### 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*]

### 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

## Quality Assurance



**Menu**

- [Definitions](#)
- [Harvest status](#)
- [Bitpreservation](#)
- [Quality Assurance](#)
- [Viewerproxy Status](#)
- [Systemstate](#)

## Viewerproxy Status

Current Viewerproxy status:

Currently does `_not_` collect missing URLs.  
 Current list of missing URLs contains 0 URLs.  
 No jobs have been chosen for viewerproxy index.

If the frame is empty, either the viewerproxy hasn't been started, or your web browser has not been configured to use it.

### Missing URL collection

- [Start collecting URLs](#)
- [Stop collecting URLs](#)
- [Clear collected URLs](#)
- [Show collected URLs](#)

### Browsing jobs in the viewerproxy

Use these pages to select the index for viewerproxy browsing:

- [Selective harvest history](#)
- [Snapshot harvest history](#)

NetarchiveSuite Version: 3.6.0 status

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

## System State

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



**Menu**

- Definitions
- Harvest status
- Bitpreservation
- Quality Assurance
- Systemstate
  - Overview of the system state

### Overview of the system state

Organisation (show all)	Machine	Port	Application	Index (show all)	
SB	kb-test-adm-001	8082	BitarchiveMonitorServer	0	Jun 13, 2008 12:21:49 PM dk.netark INFO: BatchReplyMessage: ' BatchReplyMessage for batch job ID FilesProcessed = 1 FilesFailed = 0  More...
SB	sb-test-acs-001	8075	ViewerProxy	0	Jun 13, 2008 12:21:42 PM dk.netark INFO: Starting viewerproxy jetty c
SB	sb-test-bar-001	8075	BitarchiveServer	0	Jun 13, 2008 12:21:48 PM dk.netark INFO: Finished batch job dk.netark
SB	sb-test-har-001	8081	HarvestControllerServer	0	Jun 13, 2008 12:21:38 PM dk.netark INFO: HarvestControllerServer star
SB	sb-test-har-001	8081	SideKick	0	Jun 13, 2008 12:21:37 PM dk.netark INFO: Initializing a JMS connectic


NetarchiveSuite Version: 3.6.0 status RELEASE, TEST2

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

This can be narrowed down to single "Organisations" / "Machines" / "Ports" / "Applications" 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.

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



**Menu**

- Definitions
- Harvest status
- Bitpreservation
- Quality Assurance
- Systemstate
  - Overview of the system state

### Overview of the system state

Organisation	Machine	Port	Application (show all)	Index (show all)	
KB	kb-test-har-001	8081	HarvestControllerServer	0	Jun 13, 2008 1:17:04 PM dk.netarkivet.harvest INFO: Job ID: 4, Harvest ID: 2 RUNNING timestamp discovered 2008-06-13T11:17:04Z 1102 61
KB	kb-test-har-002	8081	HarvestControllerServer	0	Jun 13, 2008 12:21:28 PM dk.netarkivet.harves INFO: HarvestControllerServer started.
KB	kb-test-har-002	8082	HarvestControllerServer	0	Jun 13, 2008 12:21:50 PM dk.netarkivet.harves INFO: HarvestControllerServer started.
SB	sb-test-har-001	8081	HarvestControllerServer	0	Jun 13, 2008 12:21:38 PM dk.netarkivet.harves INFO: HarvestControllerServer started.

NetarchiveSuite Version: 3.6.0 st

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

## Bit Preservation

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 two bitarchive instances – **SB** and **KB**.

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.

## Missing Files

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.



**Menu**

- Definitions
- Harvest status
- Bitpreservation
  - Filestatus
- Quality Assurance
- Systemstate

## Missing Files

Missing files for SB:

8-2-20080612085429-00000-kb-test-har-001.kb.dk.arc	<input checked="" type="checkbox"/> Get info
10-2-20080612085610-00000-kb-test-har-001.kb.dk.arc	<input checked="" type="checkbox"/> Get info

Status:

State	Checksum
Admin Data -	7682a44411067e15aab18ec0343f4d0d
SB	UPLOAD_COMPLETED No checksum
KB	UPLOAD_COMPLETED 7682a44411067e15aab18ec0343f4d0d

Add to archive

7-2-20080612085348-00002-kb-test-har-001.kb.dk.arc	<input checked="" type="checkbox"/> Get info
--	--

Status:

State	Checksum
Admin Data -	f601b4f165fe7dbfc358e3789f6eaabd
SB	UPLOAD_COMPLETED No checksum
KB	UPLOAD_COMPLETED f601b4f165fe7dbfc358e3789f6eaabd

Add to archive

4-metadata-1.arc  Get info

Status:

State	Checksum
Admin Data -	b66cfc990b22f5ff9d8929fdd0475915
SB	UPLOAD_COMPLETED No checksum
KB	UPLOAD_COMPLETED b66cfc990b22f5ff9d8929fdd0475915

Add to archive

Change the infobox for  files


Change  that can be added

NetarchiveSuite Version: 3.6.0 status R

If the two “remaining” checksums (On the screen dump “273356f23254684051dc99ab7c0e8937”) 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.

## Checksum Errors



**Menu**

- Definitions
- Harvest status
- Bitpreservation
  - Filestatus
- Quality Assurance
- Systemstate

## Files with checksum errors

Status for the file: 2-1-20080606122749-00001-kb-test-har-002.kb.dk.arc

Admin State	Checksum
	f255278dfa38100993bba2e1187c09cf
SB	UPLOAD_COMPLETED f255278dfa38100993bba2e1187c09cf
KB	UPLOAD_COMPLETED 545647ecb2e258c73dff0a44c9eb9961

Insert password:

Files with checksum errors in KB:

2-1-20080606122749-00001-kb-test-har-002.kb.dk.arc [Info](#)

NetarchiveSuite Version: 3.6.0 status R

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 Location XX**]. The bad file will not be completely removed but moved to an 'attic' directory on the bitarchive server holding the bad file.

edit