





Why Heritrix?

- Internet Archive's web-scale, archival-quality web crawler project
- Open-source and extensible
- Written in Java and used in CiteSeer

Download/untar/cd bin

<http://crawler.archive.org/index.html> Go to sourceforge downloads page and get version 1.14.3

```
Applications Places [Icons] [System] [Network] [Places] [Terminal] [Help]
sdas@ubuntu: ~/setups/heritrix-1.14.3/bin 1:13 PM sdas
File Edit View Search Terminal Help
sdas@ubuntu:~/setups/heritrix-1.14.3/bin$ pwd
/home/sdas/setups/heritrix-1.14.3/bin
sdas@ubuntu:~/setups/heritrix-1.14.3/bin$ echo $JAVA_HOME
/home/sdas/setups/jdk1.7.0_51
sdas@ubuntu:~/setups/heritrix-1.14.3/bin$ echo $HERITRIX_HOME
/home/sdas/setups/heritrix-1.14.3
sdas@ubuntu:~/setups/heritrix-1.14.3/bin$ ls
arcreader          extractor.cmd      heritrix.cmd      make_reports.pl
arcreader.cmd     foreground_heritrix  hoppath.pl
cmdline-jmxclient-0.10.5.jar foreground_heritrix.cmd  htmlextractor
extractor         heritrix          htmlextractor.cmd
sdas@ubuntu:~/setups/heritrix-1.14.3/bin$ ./heritrix -a gsdas:gsdas
Wed Aug 13 13:13:04 SGT 2014 Starting heritrix
Heritrix 1.14.3 is running.
Web console is at: http://127.0.0.1:8080
Web console login and password: gsdas/gsdas
sdas@ubuntu:~/setups/heritrix-1.14.3/bin$ █
```

Applications Places    

Heritrix: Login x

localhost:8080/login.jsp;jsessionid=aofsjspy76pt

HERITRIX

Login

Username:

Password:

Login

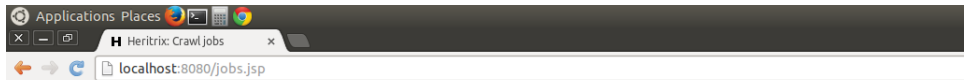
HERITRIX Status as of **Aug. 13, 2014 04:31:50 GMT** Alerts: no alerts
HOLDING JOBS
Admin Console 0 jobs pending, 1 completed

Crawler Status: **HOLDING JOBS** | [Start](#)

Jobs None running 0 pending, 1 completed Alerts: 0 (0 new)	Memory 14670 KB used 57344 KB current heap 232960 KB max heap
--	---

[Refresh](#)

[Shut down Heritrix software](#) | [Logout](#)



Status as of **Aug. 13, 2014 05:35:17 GMT** Alerts: no alerts

HOLDING JOBS

Crawl jobs

0 jobs pending, 3 completed

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

Create New Job





- [Based on existing job](#)
- [Based on a recovery](#)
- [Based on a profile](#)
- [With defaults](#)

Pending Jobs (0)

Completed Jobs(3)

UID	Name	Status	Options						
20140813044606002	mirror_crawl	Finished - Ended by operator	Crawl order	Crawl report	Seeds report	Seed file	Logs	Journal	Delete
20140813043422142	crawl_test	Finished - Ended by operator	Crawl order	Crawl report	Seeds report	Seed file	Logs	Journal	Delete
20140812043421791	default	Finished - Ended by operator	Crawl order	Crawl report	Seeds report	Seed file	Logs	Journal	Delete

Identifier: org.archive.crawler:jmxport=8849,name=Heritrix,type=CrawlService,guiport=8080,host=ubuntu

Applications Places    

Heritrix: New crawl job x

localhost:8080/jobs/new.jsp

HERITRIX Status as of **Aug. 13, 2014 04:32:09 GMT** Alerts: no alerts
HOLDING JOBS
New crawl job 0 jobs pending, 1 completed

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

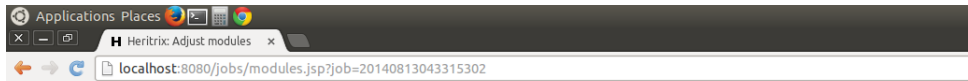
Create new crawl job based on default profile

Name of new job:

Description:

Seeds:

Identifier: org.archive.crawler:jmxport=8849,name=Heritrix,type=CrawlService,guiport=8080,host=ubuntu



Status as of **Aug. 13, 2014 04:33:15 GMT** Alerts: no alerts
HOLDING JOBS

Adjust modules 0 jobs pending, 1 completed

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

[Job crawl_test](#): [Modules](#) [Submodules](#) [Settings](#) [Overrides](#) [Refinements](#) [Submit job](#)

Select Modules and Add/Remove/Order Processors

Use this page to choose the main modules Heritrix should using crawling and to add/remove/order processors in each step of the processing chain. Go to the [Setting](#) modules and processors.

Select Crawl Scope

Current selection: org.archive.crawler.deciderules.DecidingScope
DecidingScope. A Scope that applies one or more DecideRules to determine whether a URI is accepted or rejected (returns false).

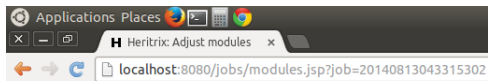
Available alternatives:

Select URI Frontier

Current selection: org.archive.crawler.frontier.BdbFrontier
BdbFrontier. A Frontier using BerkeleyDB Java Edition databases for persistence to disk.

Available alternatives:

Select Pre Processors *Processors that should run before any fetching*



Select Fetchers *Processors that fetch documents using various protocols*

org.archive.crawler.fetcher.FetchDNS	Down Remove Info
org.archive.crawler.fetcher.FetchHTTP	Up Remove Info
org.archive.crawler.prefetch.Preselector	Add

Select Extractors *Processors that extracts links from URIs*





org.archive.crawler.extractor.ExtractorHTTP	Down Remove Info
org.archive.crawler.extractor.ExtractorHTML	Up Down Remove Info
org.archive.crawler.extractor.ExtractorCSS	Up Down Remove Info
org.archive.crawler.extractor.ExtractorJS	Up Down Remove Info
org.archive.crawler.extractor.ExtractorSWF	Up Remove Info
org.archive.crawler.prefetch.Preselector	Add

Select Writers *Processors that write documents to archive files*

org.archive.crawler.writer.ARCWriterProcessor	Remove Info
org.archive.crawler.prefetch.Preselector	Add

Select Post Processors *Processors that do cleanup and feed the Frontier with new URIs*

org.archive.crawler.postprocessor.CrawlStateUpdater	Down Remove Info
org.archive.crawler.postprocessor.LinksScoper	Up Down Remove Info
org.archive.crawler.postprocessor.FrontierScheduler	Up Remove Info

Applications Places    

Heritrix: Submodules x

localhost:8080/jobs/submodules.jsp?job=20140813043422142

HERITRIX Status as of **Aug. 13, 2014 04:34:22 GMT** Alerts: no alerts
HOLDING JOBS
Submodules 0 jobs pending, 1 completed

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

[Job crawl_test](#): [Modules](#) [Submodules](#) [Settings](#) [Overrides](#) [Refinements](#) [Submit job](#)

Add/Remove/Order Submodules

Use this page to add/remove/order submodules. Go to the [Settings](#) page to complete configuration of added submodules (e.g. To add the particular regex to an added authentication information into an added RFC2617 credential).

crawl-order





scope DecidingScope. A Scope that applies one or more DecideRules to determine whether a URI is accepted or rejected (returns false).To change scope, go to the *Modules* tab.

decide-rules

rules

- rejectByDefault [Up](#) [Down](#) [Remove](#) [org.archive.crawler.deciderules.RejectDecideRule ?](#)
- acceptIfSurtPrefixed [Up](#) [Down](#) [Remove](#) [org.archive.crawler.deciderules.SurtPrefixedDecideRule ?](#)
- rejectIfTooManyHops [Up](#) [Down](#) [Remove](#) [org.archive.crawler.deciderules.TooManyHopsDecideRule ?](#)
- acceptIfTranscluded [Up](#) [Down](#) [Remove](#) [org.archive.crawler.deciderules.TransclusionDecideRule ?](#)
- rejectIfPathological [Up](#) [Down](#) [Remove](#) [org.archive.crawler.deciderules.PathologicalPathDecideRule ?](#)
- rejectIfTooManyPathSegs [Up](#) [Down](#) [Remove](#) [org.archive.crawler.deciderules.TooManyPathSegmentsDecideRule ?](#)
- acceptIfPrerequisite [Up](#) [Down](#) [Remove](#) [org.archive.crawler.deciderules.PrerequisiteAcceptDecideRule ?](#)

Name: Type: [org.archive.crawler.deciderules.AcceptDecideRule](#) [Add](#)

Applications Places    

Heritrix: Submodules x

localhost:8080/jobs/submodules.jsp?job=20140813043422142

Frontier Frontier: A frontier using BerkeleyDB Java Edition databases for persistence to disk. To change frontier, go to the *Modules* tab.

uri-canonicalization-rules

- Lowercase [Up](#) [Down](#) [Remove](#) org.archive.crawler.url.canonicalize.LowercaseRule ?
- Userinfo [Up](#) [Down](#) [Remove](#) org.archive.crawler.url.canonicalize.StripUserInfoRule ?
- WWW[0-9]* [Up](#) [Down](#) [Remove](#) org.archive.crawler.url.canonicalize.StripWWWRule ?
- SessionIDs [Up](#) [Down](#) [Remove](#) org.archive.crawler.url.canonicalize.StripSessionIDs ?
- SessionCFIDs [Up](#) [Down](#) [Remove](#) org.archive.crawler.url.canonicalize.StripSessionCFIDs ?
- QueryStringPrefix [Up](#) [Down](#) [Remove](#) org.archive.crawler.url.canonicalize.FixupQueryString ?

Name: Type: org.archive.crawler.url.canonicalize.LowercaseRule ▼

pre-fetch-processors Processors to run prior to fetching anything from the network. To change pre-fetch-processors, go to the *Modules* tab.

Preselector

Preselector#decide-rules

rules

Name: Type: org.archive.crawler.deciderules.AcceptDecideRule ▼

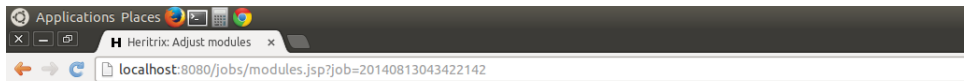
Preprocessor

Preprocessor#decide-rules

rules

Name: Type: org.archive.crawler.deciderules.AcceptDecideRule ▼

fetch-processors Processors that fetch documents. To change fetch-processors, go to the *Modules* tab.



org.archive.crawler.fetcher.FetchHTTP [Up](#) [Remove](#) [Info](#)
org.archive.crawler.prefetch.Preselector

Select Extractors *Processors that extracts links from URIs*

org.archive.crawler.extractor.ExtractorHTTP [Down](#) [Remove](#) [Info](#)
org.archive.crawler.extractor.ExtractorHTML [Up](#) [Down](#) [Remove](#) [Info](#)
org.archive.crawler.extractor.ExtractorCSS [Up](#) [Down](#) [Remove](#) [Info](#)
org.archive.crawler.extractor.ExtractorJS [Up](#) [Down](#) [Remove](#) [Info](#)
org.archive.crawler.extractor.ExtractorSWF [Up](#) [Remove](#) [Info](#)
org.archive.crawler.prefetch.Preselector

Select Writers *Processors that write documents to archive files*





org.archive.crawler.writer.ARCWriterProcessor [Remove](#) [Info](#)
org.archive.crawler.writer.MirrorWriterProcessor

Select Post Processors *Processors that do cleanup and feed the Frontier with new URIs*

org.archive.crawler.postprocessor.CrawlStateUpdater [Down](#) [Remove](#) [Info](#)
org.archive.crawler.postprocessor.LinksScoper [Up](#) [Down](#) [Remove](#) [Info](#)
org.archive.crawler.postprocessor.FrontierScheduler [Up](#) [Remove](#) [Info](#)
org.archive.crawler.prefetch.Preselector

Select Statistics Tracking

org.archive.crawler.admin.StatisticsTracker [Remove](#) [Info](#)

Applications Places    

Heritrix: Configure settin x

localhost:8080/jobs/configure.jsp?job=20140813043422142

HERITRIX

Status as of **Aug. 13, 2014 04:36:11 GMT** Alerts: no alerts
HOLDING JOBS

Configure settings 0 jobs pending, 1 completed

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

Job crawl_test: [Modules](#) [Submodules](#) [Settings](#) [Overrides](#) [Refinements](#) [Submit job](#)

[View expert settings](#)

Meta data

Description:

Crawl Operator:

Crawl Organization:

Crawl Job Recipient:

crawl-order

? Heritrix crawl order.

max-bytes-download: ?

max-document-download: ?

max-time-sec: ?

max-toe-threads: ?

scope

? DecidingScope.





enabled: ?

decide-rules

? DecideRuleSequence.

? This is a list of DecideRules to be applied in sequence.

? Select DecideRules whose names match REGEXE decision

Applications Places    

Heritrix: Configure settin x

localhost:8080/jobs/configure.jsp?job=20140813043422142

rules

rejectByDefault

? This is a list of DecideRules to be applied in sequence.

? RejectDecideRule: always gives REJECT decision.

acceptIfSurtPrefixed

? SurtPrefixedDecideRule.

decision:

? ACCEPT

surts-source-file:

?

seeds-as-surt-prefixes:

? True

rejectIfTooManyHops

? TooManyHopsDecideRule.

max-hops:

? 2

acceptIfTranscluded

? TransclusionDecideRule.

max-trans-hops:

? 3

max-speculative-hops:

? 1

rejectIfPathological

? PathologicalPathDecideRule.

max-repetitions:

? 2

rejectIfTooManyPathSegs

? TooManyPathSegmentsDecideRule.

max-path-depth:

? 20

acceptIfPrerequisite

? PrerequisiteAcceptDecideRule.

http-headers

? HTTP headers.

user-agent:

? Mozilla/5.0 (compatible; heritrix/1.14.3 +http://www.cse.psu.edu)

from:

? gsdas@cse.psu.edu

robots-honoring-policy





? Robots honoring policy

type:

? classic

masquerade:

? False

Applications Places    

Heritrix: Configure settin x

localhost:8080/jobs/configure.jsp?job=20140813043422142

rules

rejectByDefault

? This is a list of DecideRules to be applied in sequence.

? RejectDecideRule: always gives REJECT decision.

acceptIfSurtPrefixed

? SurtPrefixedDecideRule.

decision:

? ACCEPT

surts-source-file:

?

seeds-as-surt-prefixes:

? True

rejectIfTooManyHops

? TooManyHopsDecideRule.

max-hops:

? 2

acceptIfTranscluded

? TransclusionDecideRule.

max-trans-hops:

? 3

max-speculative-hops:

? 1

rejectIfPathological

? PathologicalPathDecideRule.

max-repetitions:

? 2

rejectIfTooManyPathSegs

? TooManyPathSegmentsDecideRule.

max-path-depth:

? 20

acceptIfPrerequisite

? PrerequisiteAcceptDecideRule.

http-headers

? HTTP headers.

user-agent:

? Mozilla/5.0 (compatible; heritrix/1.14.3 +http://www.cse.psu.edu)

from:

? gsdas@cse.psu.edu

robots-honoring-policy

? Robots honoring policy

type:

? classic

masquerade:

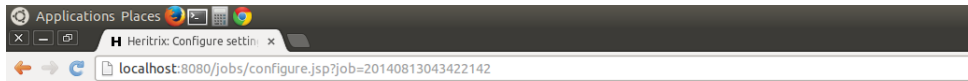
? False

frontier

- ? delay-factor:
- ? max-delay-ms:
- ? min-delay-ms:
- ? respect-crawl-delay-up-to-secs:
- ? max-retries:
- ? retry-delay-seconds:
- ? preference-embed-hops:
- ? total-bandwidth-usage-KB-sec:
- ? pause-at-start:
- ? pause-at-finish:
- ? source-tag-seeds:

pre-fetch-processors

- ? **Preselector**
Preselector.
enabled:
- ? **Preselector#decide-rules**
DecideRules which, if their final decision is REJECT, prevent this Processor from running.
- ? **Preprocessor**
Precondition enforcer
enabled:
- ? **Preprocessor#decide-rules**
DecideRules which, if their final decision is REJECT, prevent this Processor from running.



- HTTP#decide-rules** ? DecideRules which, if their final decision is REJECT, prevent this Processor from running.
- midfetch-decide-rules** ? DecideRules which, if final decision is REJECT, abort fetch after headers before all content is read.
- timeout-seconds: ?
- max-length-bytes: ?
- extract-processors** ? Processors that extract new URIs from fetched documents.
- ExtractorHTTP** ? HTTP extractor.
- enabled: ?
- ExtractorHTTP#decide-rules** ? DecideRules which, if their final decision is REJECT, prevent this Processor from running.
- ExtractorHTML** ? HTML extractor.
- enabled: ?
- ExtractorHTML#decide-rules** ? DecideRules which, if their final decision is REJECT, prevent this Processor from running.
- ExtractorCSS** ? CSS Extractor.
- enabled: ?
- ExtractorCSS#decide-rules** ? DecideRules which, if their final decision is REJECT, prevent this Processor from running.
- ExtractorJS** ? JavaScript extractor.
- enabled: ?
- ExtractorJS#decide-rules** ? DecideRules which, if their final decision is REJECT, prevent this Processor from running.
- ExtractorSWF** ? Flash extractor.
- enabled: ?
- ExtractorSWF#decide-rules** ? DecideRules which, if their final decision is REJECT, prevent this Processor from running.
- write-processors** ? Processors that write documents to archives.
- Archiver** ? ARCWriter processor
- enabled: ?



Status as of **Aug. 13, 2014 04:39:14 GMT** Alerts: no alerts

HOLDING JOBS

Crawl jobs

1 jobs pending, 1 completed

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

Job created

Create New Job

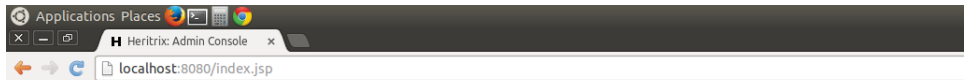
- [Based on existing job](#)
- [Based on a recovery](#)
- [Based on a profile](#)
- [With defaults](#)

Pending Jobs (1)

Name	Status	Options
crawl_test	Pending	View order Edit configuration Journal Delete

Completed Jobs(1)

UID	Name	Status	Options
20140812043421791	default	Finished - Ended by operator	Crawl order Crawl report Seeds report Seed file Logs Journal Delete



Status as of **Aug. 13, 2014 04:39:54 GMT** Alerts: no alerts

CRAWLING JOBS

PAUSED job: *crawl_test*

Admin Console

0 jobs pending, 1 completed

2 URIs in 3s (0.67/sec)

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

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

Jobs

Paused: *crawl_test*
0 pending, 1 completed

Alerts: 0 (0 new)

Memory

36329 KB used
74240 KB current heap
232960 KB max heap

Job Status: **PAUSED** | [Resume](#) | [Checkpoint](#) | [Terminate](#)

Rates

0.67 URIs/sec (0.67 avg)
0 KB/sec (0 avg)

Time

3s elapsed
1s remaining (estimated)

Totals

downloaded 2  1 queued

3 total downloaded and queued

356 B crawled (356 B novel)

Load

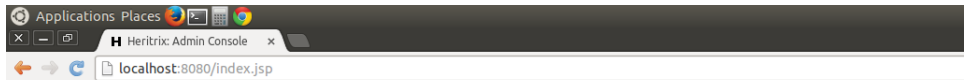
0 active of 50 threads
1 congestion ratio
1 deepest queue
1 average depth

Paused Operations

[View or Edit Frontier URIs](#)

[Refresh](#)

[Shut down Heritrix software](#) | [Logout](#)



Status as of **Aug. 13, 2014 04:39:54 GMT** Alerts: no alerts

CRAWLING JOBS

PAUSED job: *crawl_test*

Admin Console

0 jobs pending, 1 completed

2 URIs in 3s (0.67/sec)

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

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

Jobs

Paused: *crawl_test*
0 pending, 1 completed

Alerts: 0 (0 new)

Memory

36329 KB used
74240 KB current heap
232960 KB max heap

Job Status: **PAUSED** | [Resume](#) | [Checkpoint](#) | [Terminate](#)

Rates

0.67 URIs/sec (0.67 avg)
0 KB/sec (0 avg)

Time

3s elapsed
1s remaining (estimated)

Load

0 active of 50 threads
1 congestion ratio
1 deepest queue
1 average depth

Totals





downloaded 2  66% 1 queued
3 total downloaded and queued
356 B crawled (356 B novel)

Paused Operations

[View or Edit Frontier URIs](#)

[Refresh](#)

[Shut down Heritrix software](#) | [Logout](#)

Applications Places    

Heritrix: Reports x

localhost:8080/reports.jsp

HERITRIX Status as of **Aug. 13, 2014 04:40:13 GMT** Alerts: no alerts
CRAWLING JOBS PAUSED job: *crawl_test*
0 jobs pending, 1 completed 2 URIs in 3s (0.67/sec)

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

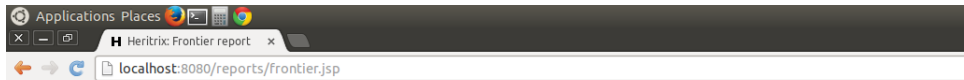
Reports on ongoing crawl/current status

- [Crawl report](#)
- [Seed report](#)
- [Frontier report](#)
1 queues: 1 active (0 in-process; 1 ready; 0 snoozed); 0 inactive; 0 retired; 0 exhausted
- [Processors report](#)
- [ToeThread report](#)
50 threads: 50 ABOUT_TO_GET_URI

The crawler generates reports when it finishes a job. Clicking here on [Force generation of end-of-crawl Reports](#) will force the writing of reports to disk. Clicking this link will force the writing of reports to disk. Clicking this link will force the writing of reports to disk. Use this facility when the crawler has hung threads that can't be interrupted.

Started crawl jobs (newest to oldest)

crawl_test *Paused*
default [Crawl report](#) [Seed report](#) *Finished - Ended by operator*



Status as of **Aug. 13, 2014 04:40:23 GMT** Alerts: no alerts

CRAWLING JOBS

PAUSED job: *crawl_test*

Frontier report

0 jobs pending, 1 completed

2 URIs in 3s (0.67/sec)

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

Frontier report - 201408130440
Job being crawled: *crawl_test*

----- STATS -----

Discovered: 3
Queued: 1
Finished: 2
Successfully: 2
Failed: 0
Disregarded: 0

----- QUEUES -----

Already included size: 3
pending: 0

All class queues map size: 1
Active queues: 1
In-process: 0
Ready: 1
Snoozed: 0
Inactive queues: 0
Retired queues: 0
Exhausted queues: 0

----- IN-PROCESS QUEUES -----

----- READY QUEUES -----

Queue *clgiles.ist.psu.edu*
1 items
last enqueued: <http://clgiles.ist.psu.edu/>

Applications Places Heritrix: View logs x
localhost:8080/logs.jsp



Status as of **Aug. 13, 2014 04:40:39 GMT** Alerts: no alerts
CRAWLING JOBS PAUSED job: *crawl_test*
0 jobs pending, 1 completed 2 URIs in 3s (0.67/sec)

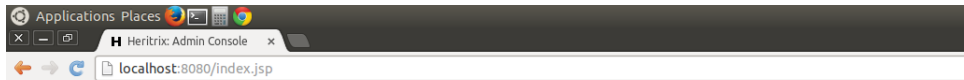
[View logs](#)
[Console](#) [Jobs](#) [Profiles](#) [Logs](#) [Reports](#) [Setup](#) [Help](#)

View: [crawl.log](#) By: [Line number](#) Refresh time:
[local-errors.log](#) [Time stamp](#)
[progress-statistics.log](#) [Regular expression](#) Lines to show:
[runtime-errors.log](#) Tail
[uri-errors.log](#)

crawl.log for crawl_test			
2014-08-13T04:39:51.318Z	1	61	dns:clgiles.ist.psu.edu P http://clgiles.ist.psu.edu/ text/dns #046 20140813043951020+
2014-08-13T04:39:53.982Z	200	47	http://clgiles.ist.psu.edu/robots.txt P http://clgiles.ist.psu.edu/ text/plain #045 20140813043953020+

[Rotate crawler logs](#)

Identifier: org.archive.crawler:jmxport=8849,name=Heritrix,type=CrawlService,guiport=8080,host=ubuntu



Status as of **Aug. 13, 2014 04:41:05 GMT** Alerts: no alerts

CRAWLING JOBS

RUNNING job: *crawl_test*

Admin Console

0 jobs pending, 1 completed

4 URIs in 9s (0.67/sec)

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

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

Jobs

Running: *crawl_test*
0 pending, 1 completed

Alerts: 0 (0 new)

Memory

37907 KB used
74240 KB current heap
232960 KB max heap

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

Rates

0.67 URIs/sec (0.67 avg)
0 KB/sec (0 avg)


Time

9s elapsed
1m14s remaining (estimated)

Load

0 active of 50 threads
1 congestion ratio
30 deepest queue
30 average depth

Totals

downloaded 4  **11%** 30 queued
34 total downloaded and queued
25 KB crawled (25 KB novel)

[Refresh](#)

[Shut down Heritrix software](#) | [Logout](#)



Status as of **Aug. 13, 2014 04:41:14 GMT** Alerts: no alerts

CRAWLING JOBS

RUNNING job: *crawl_test*

Crawl job report

0 jobs pending, 1 completed

6 URIs in 18s (0.27/sec)

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

Job name: *crawl_test* Processed docs/sec: 0.27 (0.33)

Status: Running Processed KB/sec: 2 (1)

Time: 18 sec. Total data written: 31 KB

— URIs —

Discovered: 34 ?

Queued: 28 ?

In progress: 0 ?

Total Successfully Failed Disregarded

Finished: 6 6 0 0

Status code

Documents

HTTP-200-Success-OK  5

DNS-1-OK  1

File type

Documents

Data

text/plain  2 1.4 KB





image/jpeg  1 15 KB

text/css  1 4.4 KB

text/dns  1 61 B

text/html  1 10 KB

Hosts Documents Data Time since last URI finished

Applications Places    

Heritrix: View logs x

localhost:8080/logs.jsp

HERITRIX Status as of **Aug. 13, 2014 04:42:41 GMT** Alerts: no alerts
 CRAWLING JOBS RUNNING job: *crawl_test*
 View logs 0 jobs pending, 1 completed 10 URIs in 1m45s (0/sec)

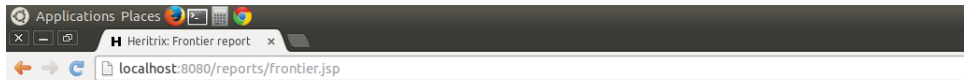
[Console](#) [Jobs](#) [Profiles](#) **[Logs](#)** [Reports](#) [Setup](#) [Help](#)

View: crawl.log By: [Line number](#) Refresh time:
[local-errors.log](#) [Time stamp](#)
[progress-statistics.log](#) [Regular expression](#) Lines to show:
[runtime-errors.log](#) Tail
[uri-errors.log](#)

crawl.log for crawl_test

2014-08-13T04:39:51.318Z	1	61	dns:clgiles.ist.psu.edu P	http://clgiles.ist.psu.edu/ text/dns #046	20140813043951020+
2014-08-13T04:39:53.982Z	200	47	http://clgiles.ist.psu.edu/robots.txt P	http://clgiles.ist.psu.edu/ text/plain #045	20140813043953982Z
2014-08-13T04:41:00.425Z	200	10232	http://clgiles.ist.psu.edu/ - - text/html #036	20140813044059768+597	shal:4JQQWIYHNBFI
2014-08-13T04:41:04.010Z	200	14927	http://clgiles.ist.psu.edu/lee-giles.jpg E	http://clgiles.ist.psu.edu/ image/jpeg #036	20140813044104010Z
2014-08-13T04:41:09.363Z	200	4234	http://clgiles.ist.psu.edu/style.css E	http://clgiles.ist.psu.edu/ text/css #049	20140813044109363Z
2014-08-13T04:41:12.335Z	200	894	http://clgiles.ist.psu.edu/favicon.ico E	http://clgiles.ist.psu.edu/ text/plain #046	20140813044112335Z
2014-08-13T04:41:31.536Z	200	410810	http://clgiles.ist.psu.edu/pubs/PLoSONE-2014.pdf L	http://clgiles.ist.psu.edu/ applica	20140813044131536Z
2014-08-13T04:41:52.147Z	200	2142	http://clgiles.ist.psu.edu/jobs.shtml L	http://clgiles.ist.psu.edu/ text/html #005	20140813044152147Z
2014-08-13T04:41:55.137Z	200	1537	http://clgiles.ist.psu.edu/pictures.shtml L	http://clgiles.ist.psu.edu/ text/html #027	20140813044155137Z
2014-08-13T04:42:40.131Z	200	1943563	http://clgiles.ist.psu.edu/papers/JCDL2009-random-forests-disambiguation.pdf L	http://	20140813044240131Z

Identifier: org.archive.crawler:jmxport=8849,name=Heritrix,type=CrawlService,guiport=8080,host=ubuntu



Status as of **Aug. 13, 2014 04:43:09 GMT** Alerts: no alerts

CRAWLING JOBS

RUNNING job: *crawl_test*

Frontier report

0 jobs pending, 1 completed

11 URIs in 2m13s (0.05/sec)

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

Frontier report - 201408130443
Job being crawled: *crawl_test*

----- STATS -----

Discovered: 40
Queued: 29
Finished: 11
Successfully: 11
Failed: 0
Disregarded: 0





----- QUEUES -----

Already included size: 40
pending: 0

All class queues map size: 1
Active queues: 1
In-process: 1
Ready: 0
Snoozed: 0
Inactive queues: 0
Retired queues: 0
Exhausted queues: 0

----- IN-PROCESS QUEUES -----

Queue *clgiles.ist.psu.edu*
29 items
last enqueued: <http://clgiles.ist.psu.edu/genealogy-giles.png>
last peeked: <http://clgiles.ist.psu.edu/bio.shtml>
total expended: 0 (total budget: -1)

Applications Places    

Heritrix: Frontier report x

localhost:8080/reports/frontier.jsp

pending: 0

All class queues map size: 1
Active queues: 1
 In-process: 1
 Ready: 0
 Snoozed: 0
Inactive queues: 0
Retired queues: 0
Exhausted queues: 0

----- IN-PROCESS QUEUES -----

Queue clgiles.ist.psu.edu
29 items
 last enqueued: http://clgiles.ist.psu.edu/genealogy-giles.png
 last peeked: http://clgiles.ist.psu.edu/bio.shtml
total expended: 0 (total budget: -1)
active balance: 3000
last(avg) cost: 0(0)

----- READY QUEUES -----

----- SNOOZED QUEUES -----

----- LONGEST QUEUE -----

Queue clgiles.ist.psu.edu
29 items
 last enqueued: http://clgiles.ist.psu.edu/genealogy-giles.png
 last peeked: http://clgiles.ist.psu.edu/bio.shtml
total expended: 0 (total budget: -1)
active balance: 3000
last(avg) cost: 0(0)

----- INACTIVE QUEUES -----

```
Applications Places [Icons] [Windows] [System Tray] 1:13 PM sdas [Settings]
sdas@ubuntu: ~/setups/heritrix-1.14.3/jobs/crawl_test-20140813043422142/arcs
File Edit View Search Terminal Help
sdas@ubuntu:~/setups/heritrix-1.14.3/jobs/crawl_test-20140813043422142$ pwd
/home/sdas/setups/heritrix-1.14.3/jobs/crawl_test-20140813043422142
sdas@ubuntu:~/setups/heritrix-1.14.3/jobs/crawl_test-20140813043422142$ ls
arcs          frontier-report.txt  order.xml          seeds-report.txt
checkpoints  hosts-report.txt    processors-report.txt  seeds.txt
crawl-manifest.txt  logs                responsecode-report.txt  state
crawl-report.txt  mimetype-report.txt  scratch            state.job
sdas@ubuntu:~/setups/heritrix-1.14.3/jobs/crawl_test-20140813043422142$ cd arcs/
sdas@ubuntu:~/setups/heritrix-1.14.3/jobs/crawl_test-20140813043422142/arcs$ ls
IAH-20140813043951-00000-ubuntu.arc.gz
sdas@ubuntu:~/setups/heritrix-1.14.3/jobs/crawl_test-20140813043422142/arcs$ █
```



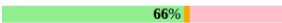
Status as of **Aug. 13, 2014 04:47:49 GMT** Alerts: no alerts
CRAWLING JOBS RUNNING job: *mirror_crawl*
Admin Console 0 jobs pending, 2 completed 2 URIs in 3s (0/sec)

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

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

Jobs Running: <i>mirror_crawl</i> 0 pending, 2 completed Alerts: 0 (0 new)	Memory 27499 KB used 88576 KB current heap 232960 KB max heap
--	---

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

Rates 0 URIs/sec (0 avg) 0 KB/sec (0 avg)	Load 0 active of 50 threads 1 congestion ratio 1 deepest queue 1 average depth
Time 3s elapsed 1s remaining (estimated)	
Totals downloaded 2  66% 1 queued 3 total downloaded and queued 356 B crawled (356 B novel)	

[Refresh](#)

[Shut down Heritrix software](#) | [Logout](#)

Applications Places    

Heritrix: View logs

localhost:8080/logs.jsp

HERITRIX Status as of **Aug. 13, 2014 04:48:55 GMT** Alerts: no alerts

CRAWLING JOBS RUNNING job: *mirror_crawl*

View logs 0 jobs pending, 2 completed 9 URIs in 1m9s (0.1/sec)

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

View: [crawl.log](#) By: [Line number](#) Refresh time:

[local-errors.log](#) [Time stamp](#) Lines to show:

[progress-statistics.log](#) [Regular expression](#)

[runtime-errors.log](#) Tail

[uri-errors.log](#)

crawl.log for mirror_crawl

2014-08-13T04:47:45.850Z	1	61	dns:clgiles.ist.psu.edu P	http://clgiles.ist.psu.edu/ text/dns #046	20140813044745834+
2014-08-13T04:47:48.460Z	200	47	http://clgiles.ist.psu.edu/robots.txt P	http://clgiles.ist.psu.edu/ text/plain #045	20140813044748460Z
2014-08-13T04:47:51.470Z	200	10232	http://clgiles.ist.psu.edu/ - - text/html #046	20140813044750841+599 sha1:4JQ0WIYHNBFI	20140813044751470Z
2014-08-13T04:47:55.054Z	200	14927	http://clgiles.ist.psu.edu/lee-giles.jpg E	http://clgiles.ist.psu.edu/ image/jpeg #045	20140813044755054Z
2014-08-13T04:48:00.391Z	200	894	http://clgiles.ist.psu.edu/favicon.ico E	http://clgiles.ist.psu.edu/ text/plain #036	20140813044800391Z
2014-08-13T04:48:03.370Z	200	4234	http://clgiles.ist.psu.edu/style.css E	http://clgiles.ist.psu.edu/ text/css #047	20140813044803370Z
2014-08-13T04:48:09.333Z	200	141613	http://clgiles.ist.psu.edu/pubs/semval2010-SEERLAB.pdf L	http://clgiles.ist.psu.edu/	20140813044809333Z
2014-08-13T04:48:27.526Z	200	9897	http://clgiles.ist.psu.edu/projects.shtml L	http://clgiles.ist.psu.edu/ text/html #026	20140813044827526Z
2014-08-13T04:48:43.476Z	404	215	http://clgiles.ist.psu.edu/pubs/CIKM2013.pdf L	http://clgiles.ist.psu.edu/ text/html #	20140813044843476Z

Identifier: org.archive.crawler:jmxport=8849,name=Heritrix,type=CrawlService,guiport=8080,host=ubuntu

```
Applications Places [Icons] [Windows] [System Tray] 1:14 PM sdas [Settings]
sdas@ubuntu: ~/setups/heritrix-1.14.3/jobs/mirror_crawl-20140813044606002/mirror/clgiles.ist.psu.edu/pubs
File Edit View Search Terminal Help
sdas@ubuntu:~/setups/heritrix-1.14.3/jobs/mirror_crawl-20140813044606002$ pwd
/home/sdas/setups/heritrix-1.14.3/jobs/mirror_crawl-20140813044606002
sdas@ubuntu:~/setups/heritrix-1.14.3/jobs/mirror_crawl-20140813044606002$ ls
checkpoints          hosts-report.txt    order.xml           seeds-report.txt
crawl-manifest.txt   logs               processors-report.txt seeds.txt
crawl-report.txt     mime-type-report.txt responsecode-report.txt state
frontier-report.txt  mirror            scratch             state.job
sdas@ubuntu:~/setups/heritrix-1.14.3/jobs/mirror_crawl-20140813044606002$ cd mirror/
sdas@ubuntu:~/setups/heritrix-1.14.3/jobs/mirror_crawl-20140813044606002/mirror$ ls
clgiles.ist.psu.edu
sdas@ubuntu:~/setups/heritrix-1.14.3/jobs/mirror_crawl-20140813044606002/mirror$ cd clgile
s.ist.psu.edu/
sdas@ubuntu:~/setups/heritrix-1.14.3/jobs/mirror_crawl-20140813044606002/mirror/clgiles.is
t.psu.edu$ ls
collaborators.shtml  index.html         projects.shtml     robots.txt
favicon.ico          lee-giles.jpg     pubs              style.css
sdas@ubuntu:~/setups/heritrix-1.14.3/jobs/mirror_crawl-20140813044606002/mirror/clgiles.is
t.psu.edu$ cd pubs/
sdas@ubuntu:~/setups/heritrix-1.14.3/jobs/mirror_crawl-20140813044606002/mirror/clgiles.is
t.psu.edu/pubs$ ls
AAAI2012-name-ethnicity.pdf  semeval2010-SEERLAB.pdf
```