
ArchiveBot Documentation

Release INF

ArchiveTeam

September 05, 2014

1 Commands	3
1.1 archive	3
1.2 abort	4
1.3 archiveonly	4
1.4 ignore	5
1.5 unignore	5
1.6 ignoreset	6
1.7 ingorereports	6
1.8 delay	6
1.9 concurrency	6
1.10 yahoo	7
1.11 expire	7
1.12 status	7
1.13 pending	7
1.14 status	8
2 Indices and tables	9

Homepage <http://www.archiveteam.org/index.php?title=ArchiveBot>

Contents:

Commands

ArchiveBot listens to commands prefixed with !.

1.1 archive

!archive URL, !a URL begin recursive retrieval from a URL

```
> !archive http://artscene.textfiles.com/litpacks/
< Archiving http://artscene.textfiles.com/litpacks/.
< Use !status 43z7a11vo6of3a7i173441dtc for updates, !abort
  43z7a11vo6of3a7i173441dtc to abort.
```

ArchiveBot does not ascend to parent links.

1.1.1 Accepted parameters

--ignore-sets SET1, ..., SETN specify sets of URL patterns to ignore:

```
> !archive http://example.blogspot.com/ncr --ignore-sets=blogs,forums
< Archiving http://example.blogspot.com/ncr.
< 14 ignore patterns loaded.
< Use !status 5sid4pgxkiu6zynhb3q1gi2s for updates, !abort
  5sid4pgxkiu6zynhb3q1gi2s to abort.
```

--no-offsite-links do not download links to offsite pages:

```
> !archive http://example.blogspot.com/ncr
>   --no-offsite-links
< Archiving http://example.blogspot.com/ncr.
< Offsite links will not be grabbed.
< Use !status 5sid4pgxkiu6zynhb3q1gi2s for updates, !abort
  5sid4pgxkiu6zynhb3q1gi2s to abort.
```

ArchiveBot's default behavior with !archive is to recursively fetch all pages that are descendants of the starting URL, as well as all linked pages and their requisites. This is often useful for preserving a page's context in time. However, this can sometimes result in an undesirably large archive. Specifying --no-offsite-links preserves recursive retrieval but does not follow links to offsite hosts.

Please note that ArchiveBot considers www.example.com and example.com to be different hosts, so if you have a website that uses both, you should not specify --no-offsite-links.

--user-agent-alias ALIAS specify a user-agent to use:

```
> !archive http://artscene.textfiles.com/litpacks/
  --user-agent-alias=firefox
< Archiving http://artscene.textfiles.com/litpacks/ .
< Using user-agent Mozilla/5.0 (Windows NT 5.1; rv:31.0)
  Gecko/20100101 Firefox/31.0.
< Use !status 43z7a11vo6of3a7i173441dtc for updates, !abort
  43z7a11vo6of3a7i173441dtc to abort.
```

This option makes the job present the given user-agent. It can be useful for archiving sites that (still) do user-agent detection. See db/user_agents for a list of recognized aliases.

--pipeline PIPELINE_ID: specify pipeline to use:

```
> !archive http://example.blogspot.com/ncr
  --pipeline=pipeline:1234567890abcdef
< Archiving http://example.blogspot.com/ncr.
< Job will run on pipeline pipeline:1234567890abcdef.
< Use !status 5sid4pgxkiu6zynhb3q1gi2s for updates, !abort
  5sid4pgxkiu6zynhb3q1gi2s to abort.
```

--phantomjs access pages via PhantomJS

--phantomjs-wait set number of seconds between PhantomJS requests; defaults to 2.0

--phantomjs-scroll maximum number of times to scroll a page in PhantomJS; defaults to 100

--no-phantomjs-smart-scroll disable PhantomJS' end-of-page detection and always scroll --phantomjs-scroll number of times; off by default

PhantomJS mode is enabled if any of the *-phantomjs* options are passed.

Tip: If you're feeling snarky or realist, you can also invoke !archive as !firstworldproblems.

1.2 abort

!abort IDENT abort a job:

```
> !abort 1q2qydhkeh3gfnrcxuf6py70b
< Initiating abort for job 1q2qydhkeh3gfnrcxuf6py70b.
```

1.3 archiveonly

!archiveonly URL, !ao URL non-recursive retrieval of the given URL:

```
> !archiveonly http://store.steampowered.com/livingroom
< Archiving http://store.steampowered.com/livingroom without
  recursion.
> Use !status 1q2qydhkeh3gfnrcxuf6py70b for updates, !abort
  1q2qydhkeh3gfnrcxuf6py70b to abort.
```

1.3.1 Accepted parameters

--ignore-sets SET1, ..., SETN specify sets of URL patterns to ignore:

```
> !archiveonly http://example.blogspot.com/ --ignore-sets=blogs,forums
< Archiving http://example.blogspot.com/ without recursion.
< 14 ignore patterns loaded.
< Use !status 5sid4pgxkiu6zynhb3q1gi2s for updates, !abort
  5sid4pgxkiu6zynhb3q1gi2s to abort.
```

--user-agent-alias ALIAS specify a user-agent to use:

```
> !archiveonly http://artscene.textfiles.com/litpacks/
  --user-agent-alias=firefox
< Archiving http://artscene.textfiles.com/litpacks/ without
  recursion.
< Using user-agent Mozilla/5.0 (Windows NT 5.1; rv:31.0)
  Gecko/20100101 Firefox/31.0.
< Use !status 43z7a11vo6of3a7i173441dtc for updates, !abort
  43z7a11vo6of3a7i173441dtc to abort.
```

This option makes the job present the given user-agent. It can be useful for archiving sites that (still) do user-agent detection. See db/user_agents for a list of recognized aliases.

--pipeline PIPELINE_ID specify pipeline to use:

```
> !archiveonly http://example.blogspot.com/
  --pipeline=pipeline:1234567890abcdef
< Archiving http://example.blogspot.com/ .
< Job will run on pipeline pipeline:1234567890abcdef.
< Use !status 5sid4pgxkiu6zynhb3q1gi2s for updates, !abort
  5sid4pgxkiu6zynhb3q1gi2s to abort.
```

--phantomjs access pages via PhantomJS

--phantomjs-wait set number of seconds between PhantomJS requests; defaults to 2.0

--phantomjs-scroll maximum number of times to scroll a page in PhantomJS; defaults to 100

--no-phantomjs-smart-scroll disable PhantomJS' end-of-page detection and always scroll --phantomjs- scroll number of times; off by default

PhantomJS mode is enabled if any of the *-phantomjs* options are passed.

1.4 ignore

!ignore IDENT PATTERN, !ig IDENT PATTERN add an ignore pattern:

```
> !ig 1q2qydhkeh3gfnrcxuf6py70b obnoxious\?foo=\d+
< Added ignore pattern obnoxious\?foo=\d+ to job
  1q2qydhkeh3gfnrcxuf6py70b.
```

The pattern must be expressed as regular expressions. For more information, see:

<http://docs.python.org/3/howto/regex.html#regex-howto> <http://docs.python.org/3/library/re.html#regular-expression-syntax>

1.5 unignore

!unignore IDENT PATTERN, !unig IDENT PATTERN remove an ignore pattern:

```
> !unig 1q2qydhkeh3gfnrcxuf6py70b obnoxious\?foo=\d+
< Removed ignore pattern obnoxious\?foo=\d+ from job
1q2qydhkeh3gfnrcxuf6py70b.
```

1.6 ignoreset

!ignoreset IDENT NAME, !igset IDENT NAME add a set of ignore patterns:

```
> !igset 1q2qydhkeh3gfnrcxuf6py70b blogs
< Added 17 ignore patterns to job 1q2qydhkeh3gfnrcxuf6py70b.
```

You may specify multiple ignore sets. Ignore sets that are unknown are, well, ignored:

```
> !igset 1q2qydhkeh3gfnrcxuf6py70b blogs, other
< Added 17 ignore patterns to job 1q2qydhkeh3gfnrcxuf6py70b.
< The following sets are unknown: other
```

Ignore set definitions can be found under db/ignore_patterns/.

1.7 ignorereports

!ignorereports IDENT on|off, !igrep IDENT on|off toggle ignore reports:

```
> !igrep 1q2qydhkeh3gfnrcxuf6py70b on
< Showing ignore pattern reports for job 1q2qydhkeh3gfnrcxuf6py70b.

> !igrep 1q2qydhkeh3gfnrcxuf6py70b off
< Suppressing ignore pattern reports for job
1q2qydhkeh3gfnrcxuf6py70b.
```

Some jobs generate ignore patterns at high speed. For these jobs, turning off ignore pattern reports may improve both the usefulness of the dashboard job log and the speed of the job.

This command is aliased as `!ligoff IDENT` and `!igon IDENT`. `!ligoff` suppresses reports; `!igon` shows reports.

1.8 delay

!delay IDENT MIN MAX, !d IDENT MIN MAX set inter-request delay:

```
> !delay 1q2qydhkeh3gfnrcxuf6py70b 500 750
< Inter-request delay for job 1q2qydhkeh3gfnrcxuf6py70b set to [500,
750 ms].
```

Delays may be any non-negative number, and are interpreted as milliseconds. The default inter-request delay range is [250, 375] ms.

1.9 concurrency

!concurrency IDENT LEVEL, !con IDENT LEVEL set concurrency level:

```
> !concurrency 1q2qydhkeh3gfnrcxuf6py70b 8
< Job 1q2qydhkeh3gfnrcxuf6py70b set to use 8 workers.
```

Adding additional workers may speed up grabs if the target site has capacity to spare, but it also puts additional pressure on the target. Use wisely.

1.10 yahoo

!yahoo IDENT set zero second delays, crank concurrency to 11:

```
> !yahoo 1q2qydhkeh3gfnrcxuf6py70b
< Inter-request delay for job 1q2qydhkeh3gfnrcxuf6py70b set to
[0, 0] ms.
< Job 1q2qydhkeh3gfnrcxuf6py70b set to use 11 workers.
```

Only recommended for use when archiving data from hosts with gobs of bandwidth and processing power (e.g. Yahoo, Google, Amazon). Keep in mind that this is likely to trigger any rate limiters that the target may have.

1.11 expire

!expire IDENT for expiring jobs, expire a job immediately:

```
> !expire 1q2qydhkeh3gfnrcxuf6py70b
< Job 1q2qydhkeh3gfnrcxuf6py70b expired.
```

In rare cases, the 48 hour timeout enforced by ArchiveBot on archive jobs is too long. This command permits faster snapshotting. It should be used sparingly; abuse is very easy to spot.

If a job's expiry timer has not yet started, this command does not affect the given job:

```
> !expire 5sid4pgxkiu6zynhbt3q1gi2s
< Job 5sid4pgxkiu6zynhbt3q1gi2s does not yet have an expiry timer.
```

This is intended to prevent expiration of active jobs.

1.12 status

!status print job summary:

```
> !status
< Job status: 0 completed, 0 aborted, 0 in progress, 0 pending
```

1.13 pending

!pending send pending queue in private message:

```
> !pending
< [privmsg] 2 pending jobs:
< [privmsg] 1. http://artscene.textfiles.com/litpacks/
(43z7a11vo6of3a7i173441dtc)
```

```
< [privmsg] 2. http://example.blogspot.com/ncr  
(5sid4pgxkiu6zynhb3q1gi2s)
```

Jobs are listed in the order that they'll be worked on. This command lists only the global queue; it doesn't yet show the status of any pipeline-specific queues.

1.14 status

!status IDENT, !status URL print information about a job or URL

For an unknown job:

```
> !status 1q2qydhkeh3gfnrcxuf6py70b  
< Sorry, I don't know anything about job 1q2qydhkeh3gfnrcxuf6py70b.
```

For a URL that hasn't been archived:

```
> !status http://artscene.textfiles.com/litpacks/  
< http://artscene.textfiles.com/litpacks/ has not been archived.
```

For a URL that hasn't been archived, but has children that have been processed before (either successfully or unsuccessfully):

```
> !status http://artscene.textfiles.com/  
< http://artscene.textfiles.com/ has not been archived.  
< However, there have been 5 download attempts on child URLs.  
< More info: http://www.example.com/#/prefixes/http://artscene.textfiles.com/
```

For an ident or URL that's in progress:

```
> !status 43z7a11vo6of3a7i173441dtc  
<  
< Downloaded 10.01 MB, 2 errors encountered  
< More info at my dashboard: http://www.example.com
```

For an ident or URL that has been successfully archived within the past 48 hours:

```
> !status 43z7a11vo6of3a7i173441dtc  
< Archived to http://www.example.com/site.warc.gz  
< Eligible for rearchival in 30h 25m 07s
```

For an ident or URL identifying a job that was aborted:

```
> !status 43z7a11vo6of3a7i173441dtc  
< Job aborted  
< Eligible for rearchival in 00h 00m 45s
```

Indices and tables

- *genindex*
- *modindex*
- *search*