

NUTCH CHEAT SHEET

What is NUTCH you ask?

[Nutch](#) is a very popular open source JAVA based search engine built on top of [Lucene](#) which is translated to C, C++, C#, Python, Perl and Ruby. It provides all of the tools you need to run your very own search engine. Current version of nutch (as of October 2010) is: 1.2

Some important gotchas on nutch:

- Founded in 2003 by Doug Cutting, the Lucene creator, and Mike Cafarella
- Apache project since 2004 (sub-project of Lucene)
- Spin-offs:
 - [Map-Reduce](#) and distributed FS → [Hadoop](#)
 - Content type detection and parsing → [Tika](#)
- Many installations in operation, mostly vertical search
- Collections typically 1 mln - 200 mln documents

Usage:

nutch [-core] COMMAND

where COMMAND is one of the commands from following table or CLASSNAME the class named CLASSNAME most commands print help when invoked w/o parameters.

Expert: -core option is for developers only. It avoids building the job jar, instead it simply includes classes compiled with ant compile-core. [NOTE: this works only for jobs executed in 'local' mode]

Command	Description	Usage	Example
crawl	one-step crawler for intranets	Crawl <urlDir> [-dir d] [-threads n] [-depth i] [-topN N] [-solr solrURL]	bin/nutch crawl urls -dir crawl -depth 3 -topN 50
readdb	read / dump crawl db	CrawlDbReader <crawldb> (-stats -dump <out_dir> -topN <nnnn> <out_dir> [<min>] -url <url>) <crawldb> directory name where crawldb is located -stats [-sort] print overall statistics to System.out [-sort] list status sorted by host -dump <out_dir> [-format normal csv] dump the whole db to a text file in <out_dir> [-format csv] dump in Csv format	1. bin/nutch readdb crawl/crawl db -dump reports/dg 2. bin/nutch readdb crawl/crawl db -stats

		<p>[-format normal] dump in standard format (default option)</p> <p>-url <url> print information on <url> to System.out</p> <p>-topN <nnnn> <out_dir> [<min>] dump top <nnnn> urls sorted by score to <out_dir></p> <p>[<min>] skip records with scores below this value. This can significantly improve performance.</p>	
convdb	convert crawl db from pre-0.9 format	<p>CrawlDbConverter <oldDb> <newDb> [-withMetadata]</p> <p>oldDb name of the crawl db that uses UTF8 class.</p> <p>newDb name of the output crawl db that will use Text class.</p> <p>withMetadata convert also all metadata keys that use UTF8 to Text.</p>	bin/nutch convdb crawl/crawl db
mergedb	merge crawl db-s, with optional filtering	<p>CrawlDbMerger <output_crawl db> <crawl db1> [<crawl db2> <crawl db3> ...] [-normalize] [-filter]</p> <p>output_crawl db output CrawlDb</p> <p>crawl db1 ... input CrawlDb-s (single input CrawlDb is ok)</p> <p>-normalize use URLNormalizer on urls in the crawl db(s) (usually not needed)</p> <p>-filter use URLFilters on urls in the crawl db(s)</p>	bin/nutch mergesegs crawl/segments -dir crawl/segments
readlink db	read / dump link db	<p>LinkDbReader <linkdb> {-dump <out_dir> -url <url>}</p> <p>-dump <out_dir> dump whole link db to a text file in <out_dir></p> <p>-url <url> print information about <url> to System.out</p>	bin/nutch readlinkdb crawl/linkdb -dump reports/link
inject	inject new urls into the database	<p>Injector <crawl db> <url_dir></p>	bin/nutch inject crawl/crawl db seed/dynamicguy/url
generate	generate new segments to fetch from crawl db	<p>Generator <crawl db> <segments_dir> [-force] [-topN N] [-numFetchers numFetchers] [-adddays numDays] [-noFilter] [-noNorm][--maxNumSegments num]</p>	<p>1. bin/nutch generate crawl/crawl db crawl/segments</p> <p>2. bin/nutch generate crawl/crawl db crawl/segment</p>

			ents -topN 100
freegen	generate new segments to fetch from text files	<p>FreeGenerator <inputDir> <segmentsDir> [-filter] [-normalize]</p> <p>inputDir input directory containing one or more input files. Each text file contains a list of URLs, one URL per line</p> <p>segmentsDir output directory, where new segment will be created</p> <p>-filter run current URLFilters on input URLs</p> <p>-normalize run current URLNormalizers on input URLs</p>	bin/nutch freegen -dir seed/dynam icguy/ crawl/segm ents/*
fetch	fetch a segment's pages	<p>Fetcher <segment> [-threads n] [-noParsing]</p>	bin/nutch fetch crawl/segm ents/20100 929225609
parse	parse a segment's pages	<p>ParseSegment segment</p>	bin/nutch parse crawl/segm ents/20100 930105013/
readseg	read / dump segment data	<p>SegmentReader (-dump ... -list ... -get ...) [general options]</p> <p>* General options:</p> <ul style="list-style-type: none"> -nocontent ignore content directory -nofetch ignore crawl_fetch directory -nogenerate ignore crawl_generate directory -noparse ignore crawl_parse directory -noparsedata ignore parse_data directory -noparsetext ignore parse_text directory <p>* SegmentReader -dump <segment_dir> <output> [general options]</p> <p>Dumps content of a <segment_dir> as a text file to <output>.</p> <p><segment_dir> name of the segment directory.</p> <p><output> name of the (non-existent) output directory.</p> <p>* SegmentReader -list (<segment_dir1> ... -dir <segments>) [general options]</p> <p>List a synopsis of segments in specified directories, or all segments in</p>	bin/nutch readseg -list crawl/segm ents/20100 930105013/

		<p>a directory <segments>, and print it on System.out</p> <p><segment_dir1> ... list of segment directories to process</p> <p>-dir <segments> directory that contains multiple segments</p> <p>* SegmentReader -get <segment_dir> <keyValue> [general options]</p> <p>Get a specified record from a segment, and print it on System.out.</p> <p><segment_dir> name of the segment directory.</p> <p><keyValue> value of the key (url).</p> <p>Note: put double-quotes around strings with spaces.</p>	
mergesegs	merge several segments, with optional filtering and slicing	<p>SegmentMerger output_dir (-dir segments seg1 seg2 ...) [-filter] [-slice NNNN]</p> <p>output_dir name of the parent dir for output segment slice(s)</p> <p>-dir segments parent dir containing several segments</p> <p>seg1 seg2 ... list of segment dirs</p> <p>-filter filter out URL-s prohibited by current URLFilters</p> <p>-slice NNNN create many output segments, each containing NNNN URLs</p>	bin/nutch mergesegs crawl/segments -i -ds
updatedb	update crawl db from segments after fetching	<p>CrawlDb <crawldb> (-dir <segments> <seg1> <seg2> ...) [-force] [-normalize] [-filter] [-noAdditions]</p> <p>crawldb CrawlDb to update</p> <p>-dir segments parent directory containing all segments to update from</p> <p>seg1 seg2 ... list of segment names to update from</p> <p>-force force update even if CrawlDb appears to be locked (CAUTION advised)</p> <p>-normalize use URLNormalizer on urls in CrawlDb and segment (usually not needed)</p> <p>-filter use URLFilters on urls in CrawlDb and segment</p> <p>-noAdditions only update already existing URLs, don't add any newly discovered URLs</p>	bin/nutch updatedb crawl/crawldb \$s2
invertlinks	create a linkdb from parsed segments	<p>CrawlDb <crawldb> (-dir <segments> <seg1> <seg2> ...) [-force] [-normalize] [-filter] [-noAdditions]</p> <p>crawldb CrawlDb to update</p> <p>-dir segments parent directory containing all segments to update from</p>	bin/nutch invertlinks crawl/linkdb crawl/segments/20100

		seg1 seg2 ... list of segment names to update from -force force update even if CrawlDb appears to be locked (CAUTION advised) -normalize use URLNormalizer on urls in CrawlDb and segment (usually not needed) -filter use URLFilters on urls in CrawlDb and segment -noAdditions only update already existing URLs, don't add any newly discovered URLs	930105013
mergelinkdb	merge linkdb-s, with optional filtering	LinkDbMerger <output_linkdb> <linkdb1> [<linkdb2> <linkdb3> ...] [-normalize] [-filter] output_linkdb output LinkDb linkdb1 ... input LinkDb-s (single input LinkDb is ok) -normalize use URLNormalizer on both fromUrls and toUrls in linkdb(s) (usually not needed) -filter use URLFilters on both fromUrls and toUrls in linkdb(s)	bin/nutch mergelinkdb crawl/indexes/ crawl/crawldb/ crawl/linkdb/ crawl/segments/20100930105013/
index	run the indexer on parsed segments and linkdb	Indexer <index> <crawldb> <linkdb> <segment> ...	bin/nutch index crawl/indexes crawl/crawldb crawl/linkdb crawl/segments/*
solrindex	run the solr indexer on parsed segments and linkdb	SolrIndexer <solr url> <crawldb> <linkdb> <segment> ...	bin/nutch solrindex http://127.0.0.1:8080/solr/ crawl/crawldb crawl/linkdb crawl/segments/*
merge	merge several segment indexes	IndexMerger [-workingdir <workingdir>] outputIndex indexesDir...	bin/nutch solrdedup http://localhost:8983/solr/

dedup	remove duplicates from a set of segment indexes	DeleteDuplicates <indexes> ...	bin/nutch dedup crawl/indexes
solrdedup	remove duplicates from solr	SolrDeleteDuplicates <solr url>	bin/nutch solrdedup http://localhost:8983/solr
plugin	load a plugin and run one of its classes main()	PluginRepository pluginId className [arg1 arg2 ...]	in/nutch plugin parse-html org.apache.nutch.parse.html.HtmlParser path/to/file.html
server	run a search server	DistributedSearch\$Server <port> <crawl dir>	bin/nutch server 1234 crawl

Nutch Documents:

Field	Stored	Indexed	analyzed
Url	YES	YES	YES
anchor	NO	YES	YES
content	NO	YES	YES
Site	YES	YES	NO
Lang	YES	YES	NO