

# Searching Web Archive Collections

Michael Stack  
Internet Archive  
The Presidio of San Francisco  
116 Sheridan Ave. San Francisco, CA 94129  
stack@archive.org

## Abstract

*Web archive collection search presents the usual set of technical difficulties searching large collections of documents. It also introduces new challenges often at odds with typical search engine usage. This paper outlines the challenges and describes adaptation of an open source search engine, Nutch, to Web archive collection search. Statistics and observations indexing and searching small to medium-sized collections are presented. We close with a sketch of how we intend to tackle the main limitation, scaling archive collection search above the current ceiling of approximately 100 million documents.*

## 1. Introduction

The Internet Archive (IA)(www.archive.org) is a 501(c)(3) non-profit organization whose mission is to build a public Internet digital library [1]. Since 1996, the IA has been busy establishing the largest public Web archive to date, hosting over 600 terabytes of data. Currently the only public access to the Web archive has been by way of the IA Wayback Machine (WM) [2] in which users enter an URL and the WM displays a list of all instances of the URL archived, distinguished by crawl date. Selecting any date begins browsing a site as it appeared then, and continued navigation pulls up nearest-matches for linked pages. The WM suffers one major shortcoming: unless you know beforehand the exact URL of the page you want to browse, you will not be able to directly access archived content. Current Web URLs and published references to historic URLs may suggest starting points, but offer little help for thorough or serendipitous exploration of archived sites. URL-only retrieval also frustrates users who are accustomed to exhaustive Google-style full text search from a simple query box. What is missing is a full text search tool that works over archived content, to better guide users 'wayback' in time.

Nutch [4] was selected as the search engine platform on which to develop Web Archive Collection (WAC) search. "Nutch is a complete open-source Web search engine package that aims to index the World Wide Web as effectively as commercial search services" [5].

Technically, Nutch provides basic search engine capability, is extensible, aims to be cost-effective, and is demonstrated capable of indexing up to 100 million documents with a convincing development story for how to scale up to billions [9].

This paper begins with a listing of challenges searching WACs. This is followed by an overview of Nutch operation to aid understanding of the next section, a description of Nutchwax, the open-source Nutch extensions made to support WAC search. Statistics on indexing rates, index sizes, and hardware are presented as well as observations on the general WAC indexing and search operation. We conclude with a sketch of how we intend to scale up to index collections of billions of documents.

## 2. Challenges Searching WACs

WACs tend to be large. A WAC usually is an aggregate of multiple, related focused Web crawls run over a distinct time period. For example, one WAC, made by the IA comprises 140 million URLs collected over 34 weekly crawls of sites pertaining to the United States 2004 Presidential election. Another WAC is the complete IA repository of more than 60 billion URLs. (Although this number includes exact or near duplicates, the largest live-Web search engine, Google, only claims to be "[s]earching 8,058,044,651 web pages" as of this writing.) WACs are also large because archives tend not to truncate Web downloads and to fetch all resources including images and streams, not just text-only resources.

A single URL may appear multiple times in a WAC. Each instance may differ radically, minimally or not at all across crawls, but all instances are referred to using the same URL. Multiple versions complicate search query and result display: Do we display all versions in search results? If not, how do we get at each instance in the collection? Do we suppress duplicates? Or do we display the latest with a count of known instances in a corner of the search result summary?

A WAC search engine gets no help from the Web-at-large serving search results. What we mean by this is that for WAC searching, after a user clicks on a search result hit, there is still work to be done. The search result must

refer the user to a viewer or replay utility – a tool like the IA WM – that knows how to fetch the found page from the WAC repository and display it as faithfully as possible. (Since this redisplay is from a server other than the page's original home, on-the-fly content rewriting is often required.) While outside of the purview of collection search, WAC tools that can reassemble the pages of the past are a critical component in any WAC search system.

### 3. Overview Of Nutch Operation

The Nutch search engine indexing process runs in a stepped, batch mode. With notable exceptions discussed later, the intent is that each step in the process can be "segmented" and distributed across machines so no single operation overwhelms as the collection grows. Also, where a particular step fails (machine crash or operator misconfiguration), that step can be restarted. A custom database, the Nutch webdb, maintains state between processing steps and across segments. An assortment of parse-time, index-time, and query-time plugins allows amendment of each processing step.

After initial setup and configuration, an operator manually steps through the following cycle indexing:

1. Ask the Nutch webdb to generate a number of URLs to fetch. The generated list is written to a "segment" directory.
2. Run the built-in Nutch fetcher. During download, an md5 hash of the document content is calculated and parsers extract searchable text. All is saved to the segment directory.
3. Update the Nutch webdb with vitals on URLs fetched. An internal database analysis step computes all in-link anchor text per URL. When finished, the results of the database in-link anchor text analysis are fed back to the segment. Cycle steps 1-3 writing new segments per new URL list until sufficient content has been obtained.
4. Index each segment's extracted page text and in-link anchor text. Index is written into the segment directory.
5. Optionally remove duplicate pages from the index.
6. Optionally merge all segment indices (Unless the index is large and needs to be distributed).

Steps 2 and 4 may be distributed across multiple machines and run in parallel if multiple segments. Steps 1, 3, and 5 require single process exclusive access to the webdb. Steps 3 and 6 require that a single process have exclusive access to all segment data. A step must complete before the next can begin.

To query, start the Nutch search Web application. Run multiple instances of the search Web application to distribute query processing. The queried server distributes the query by remotely invoking queries against all query cluster participants. (Each query cluster participant is responsible for some subset of all segments.) Queries are run against Nutch indices and return ranked Google-like search results that include

snippets of text from the pertinent page pulled from the segment-extracted text.

## 4. Nutch Adaptation

Upon consideration, WAC search needs to support two distinct modes of operations. First, WAC search should function as a Google-like search engine. In this mode, users are not interested in search results polluted by multiple duplicate versions of a single page. Phase one of the Nutch adaptation focused on this mode of operation.

A second mode becomes important when users want to study how pages change over time. Here support for queries of the form, "return all archive versions crawled in 1999 sorted by crawl date" is needed. (Satisfying queries of this specific type is what the IA WM does using a sorted flat file index to map URL and date to resource location.) Phase two added features that allow version- and date-aware querying. (All WAC plugin extensions, documentation, and scripts are open source hosted at Sourceforge under the Nutchwax project [8].)

### 4.1. Phase one

Because the WAC content already exists, previously harvested by other means, the Nutch fetcher step had to be recast to pull content from a WAC repository rather than from the live Web. At IA, harvested content is stored in the ARC file format [6]; composite log files each with many collected URLs. For the IA, an ARC-to-segment tool was written to feed ARCs to Nutch parsers and segment content writers. (Adaptation for formats other than IA ARC should be trivial.) Upon completion of phase one, using indices purged of exact duplicates, it was possible to deploy a basic WAC search that used the IA WM as the WAC viewer application.

### 4.2. Phase two

To support explicit date and date range querying using the IA 14-digit YYYYDDMMHSS timestamp format, an alternate *date* query operator implementation replaced the native Nutch YYYYMMDD format. To support retrieval of WAC documents by IA WM-like viewer applications, location information -- *collection*, *arcname* and *arcoffset* -- was added to search results as well as an operator to support exact, as opposed to fuzzy, URL querying (*exacturl*). Nutch was modified to support sorting on arbitrary fields and deduplication at query time (*sort*, *reverse*, *dedupField*, *hitsPerDup*). Here is the complete list of new query operators:

- *sort*: Field to sort results on. Default is no sort.
- *reverse*: Set to true to reverse sort. Default is false.
- *dedupField*: Field to deduplicate on. Default is 'site'.

- *hitsPerDup*: Count of *dedupField* matches to show in search results. Default 2.
- *date*: IA 14-digit timestamps. Ranges specified with '-' delimiter between upper and lower bounds.
- *arcname*: Name of ARC file that containing result found.
- *arcoffset*: Offset into *arcname* at which result begins.
- *collection*: The collection the search result belongs to.
- *exacturl*: Query for an explicit url.

Natively Nutch passes all content for which there is no explicit parser to the text/html parser. Indexing, logs are filled with skip messages from the text/html parser as it passes over audio/\*, video/\*, and image/\* content. Skipped resources get no mention in the index and so are not searchable. An alternate parser-default plugin was created to add at least base metadata of crawl date, arcname, arcoffset, type, and URL. This allows viewer applications, which need to render archived pages that contain images, stylesheets, or audio to ask of the Nutch index the location of embedded resources. Finally, an option was added to return results as XML (RSS) [7]. Upon completion of phase two, both modes of operation were possible using a single non-duplicated index.

## 5. Indexing Stats

Discussed below are details indexing two WACs: One small, the other medium-sized. All processing was done on machines of the following profile: single processor 2.80GHz Pentium 4s with 1GB of RAM and 4x400GB IDE disks running Debian GNU/Linux. Indexing, this hardware was CPU-bound with light I/O loading. RAM seemed sufficient (no swapping). All source ARC data was NFS mounted. Only documents of type text/\* or application/\* and HTTP status code 200 were indexed.

### 5.1. Small Collection

This small collection was comprised of three crawls. Indexing steps were run in series on one machine using a single disk. The collection comprised 206 ARC files, 37.2GB of uncompressed data. 1.07 million of the collection total of 1.27 million documents was indexed.

**Table 1: MIME Types**

MIME Type	Size (MB)	% Size	Incidence
text	25767.67	79.32%	1052103
text/html	22003.55	67.73%	1044250
application	6719.92	20.68%	20969
application/pdf	4837.89	14.89%	16201
application/msword	487.89	1.50%	3306

**Table 2: Timings**

Segment	Database	Index	Dedup	Merge
16h32m	2h26m	18h44m	0h01m	02h35m

Indexing took 40.3 hours to complete. The merged index size was 1.1GB, about 3% the size of the source

collection. The index plus the cleaned-up segment data – cleaning involved removal of the (recalculable) segment-level indices made redundant by the index merge – occupied 1.1GB + 4.9GB, or about 16% the size of the source collection. Uncleaned segments plus index made up about 40% the size of the source collection.

### 5.2 Medium-sized Collection

The collection was made up of 1054 ARCs, 147.2GB of uncompressed data. 4.1 million documents were indexed. Two machines were used to do the segmenting step. Subsequent steps were all run in series on a single machine using a single disk.

**Table 3: MIME Types**

MIME Type	Size (MB)	% Size	Incidence
Text	96882.84	65.32%	3974008
text/html	90319.81	60.81%	3929737
Application	50338.40	34.68%	122174
application/pdf	21320.83	14.40%	45427
application/msword	1000.70	0.67%	5468

**Table 4: Timings**

Segment	Database	Index	Dedup	Merge
12h32m 19h18m	7h23m	55h07m	0h06m	0h31m

Indexing took 99 hours of processing time (or 86.4 hours of elapsed time because segmenting was split and run concurrently on two machines). The merged index size was 5.2GB, about 4% the size of source collection. Index plus the cleaned-up segment data occupied 5.2GB + 14.5GB, or about 13.5% the size of the source collection. (Uncleaned segments plus index occupied about 22% the size of the source collection.)

## 6. Observations

Indexing big collections is a long-running manual process that currently requires intervention at each step moving the process along. Attention required compounds the more distributed the indexing is made. An early indexing of a collection of approximately 85 million documents took more than a week to complete with segmenting and indexing spread across 4 machines. Steps had to be restarted as disks overfilled and segments had to be redistributed. Little science was applied so the load was suboptimally distributed with synchronizations waiting on laggard processes. (Others close in to the Nutch project have reported similar experiences [12].) An automated means of efficiently distributing the parsing, update, and indexing work across a cluster needs to be developed. In the way of any such development are at least the following obstacles:

- Some indexing steps are currently single process.

- As the collection grows, with it grows the central webdb of page and link content. Eventually it will grow larger than any available single disk.

We estimate that with the toolset as is, given a vigilant operator and a week of time plus 4 to 5 machines with lots of disk, indexing WACs of about 100 million documents is at the limit of what is currently practical.

Adding to a page its inlink anchor-text when indexing improves search result quality. Early indexing experiments were made without the benefit of the Nutch link database — our custom fetcher step failed to properly provide link text for Nutch to exploit. Results were rich in query terms but were not what was 'expected'. A subsequent fix made link-text begin to count. Thereafter, search result quality improved dramatically.

The distributed Nutch query clustering works well in our experience, at least for low rates of access: ~1 query per second. (Search access-rates are expected to be lower for WACs than live-Web search engines.) But caches kept in the search frontend to speed querying will turn problematic with regular usage. The base Nutch (Lucene) query implementation uses one byte per document per field indexed. Additions made to support query-time deduplication and sorting share a cache that stores each search result's document URL. Such a cache of (Java) UTF-16 Java strings gets large fast. An alternate smaller memory-footprint implementation needs to be developed.

## 7. Future Work

From inception, the Nutch project has set its sights on operating at the scale of the public web and has been making steady progress addressing the difficult technical issues scaling up indexing and search. The Nutch Distributed File System (NDFS) is modeled on a subset of the Google File System (GFS) [11] and is "...a set of software for storing very large stream-oriented files over a set of commodity computers. Files are replicated across machines for safety, and load is balanced fairly across the machine set" [12]. The intent is to use NDFS as underpinnings for a distributed webdb. (It could also be used storing very large segments.) While NDFS addresses the problem of how to manage large files in a fault-tolerant way, it does not help with the even distribution of processing tasks across a search cluster. To this end, the Nutch project is working on a version of another Google innovation, MapReduce [13], "a platform on which to build scalable computing" [9]. In synopsis, if you can cast the task you wish to run on a cluster into the MapReduce mold -- think of the Python *map* function followed by *reduce* function -- then the MapReduce platform will manage the distribution of your task across the cluster in a fault-tolerant way. Mid 2005, core developers of the Nutch project are writing a Java version

of the MapReduce platform to use in a reimplementaion of Nutch as MapReduce tasks [9]. MapReduce and NDFS combined should make Nutch capable of scaling its indexing step to billions of documents.

The IA is moving its collections to the Petabox platform; racks of low power, high storage density, inexpensive rack-mounted computers [14]. The future of WAC search development will be harnessing Nutch MapReduce/NDFS development on Petabox.

## 8. Acknowledgements

Doug Cutting and all members of the IA Web Team: Michele Kimpton, Gordon Mohr, Igor Ranitovic, Brad Tofel, Dan Avery, and Karl Thiessen. The International Internet Preservation Consortium (IIPC) [3] supported the development of Nutchwax.

## 9. References

- [1] Internet Archive <http://www.archive.org>
- [2] Wayback Machine <http://www.archive.org/web/web.php>
- [3] International Internet Preservation Consortium <http://netpreserve.org>
- [4] Nutch <http://lucene.apache.org/nutch/>
- [5] Nutch: A Flexible and Scalable Open-Source Web Search Engine <http://labs.commerce.net/wiki/images/0/06/CN-TR-04-04.pdf>
- [6] ARC File Format <http://www.archive.org/web/researcher/ArcFileFormat.php>
- [7] A9 Open Search <http://opensearch.a9.com/>
- [8] Nutchwax <http://access.archive.org/projects/nutch>
- [9] MapReduce in Nutch, 20 June 2005, Yahoo!, Sunnyvale, CA, USA <http://wiki.apache.org/nutch-data/attachments/Presentations/attachments/mapred.pdf>
- [10] The Nutch Distributed File System by Michael Cafarella <http://wiki.apache.org/nutch/NutchDistributedFileSystem>
- [11] Google File System <http://www.google.com/url?sa=U&start=1&q=http://labs.google.com/papers/gfs-sosp2003.pdf&e=747>
- [12] "[nutch-dev] Experience with a big index" by Michael Cafarella <http://www.mail-archive.com/nutch-developers@lists.sourceforge.net/msg02602.html>
- [13] MapReduce: Simplified Data Processing on Large Clusters by Jeffrey Dean and Sanjay Ghemawat <http://labs.google.com/papers/mapreduce-osdi04.pdf>
- [14] Petabox <http://www.archive.org/web/petabox.php>