

Nutch - The Java Search Engine

By Tyrell Perera, Virtusa Corp. (<http://www.virtusa.com>)

1 Introduction

1.1 What is Nutch?

Nutch is an effort to build a Free and Open Source search engine. It uses Lucene for the search and index component. The fetcher (robot) has been written from scratch solely for this project.

Nutch has a highly modular architecture allowing developers to create plug-ins for activities such as media-type parsing, data retrieval, querying and clustering.

Doug Cutting is the lead developer of Nutch.

1.2 What is Lucene?

Lucene is a Free and Open Source search and index API released by the Apache Software Foundation. It is written in Java and is released under the Apache Software License.

Lucene is just the core of a search engine. As such, it does not include things like a web spider or parsers for different document formats. Instead these things need to be added by a developer who uses Lucene.

Lucene does not care about the source of the data, its format, or even its language, as long as you can convert it to text. This means you can use Lucene to index and search data stored in files: web pages on remote web servers, documents stored in local file systems, simple text files, Microsoft Word documents, HTML or PDF files, or any other format from which you can extract textual information.

Lucene has been ported or is in the process of being ported to various programming languages other than Java:

- `Lucene4c` - C
- `CLucene` - C++
- `MUTIS` - Delphi
- `NLucene` - .NET
- `DotLucene` - .NET
- `Plucene` - Perl
- `Pylucene` - Python
- `Ferret` and `RubyLucene` - Ruby

1.3 What License?

Both Nutch and Lucene are Apache projects and carry the Apache license (<http://www.opensource.org/licenses/apache2.0.php>).

2 The Design of Nutch

2.1 Core Components of Nutch

The Nutch search engine consists, very roughly, of three components:

1. The Crawler, which discovers and retrieves web pages
2. The 'WebDB', a custom database that stores known URLs and fetched page contents
3. The 'Indexer', which dissects pages and builds keyword-based indexes from them

💡 After the initial creation of an Index, it is usual to perform periodic updates of the index, in order to keep it up-to-date. We will look into the details of index maintenance in the parts following this.

2.2 The Nutch Web Application

Apart from the above three components, it has a Search Web Application. This application is a JSP application that can be configured and deployed in a servlet container.

2.3 The Nutch API

All components listed above use the nutch API. The users can utilize the API via two approaches, which depends on the task at hand.

1. Through the nutch Shell Script for administrative tasks, such as creating and maintaining indexes
2. Through the Search Web Application, in order to perform a search using keywords

The *sequence diagram* below shows how each of these components interact in implementing a Nutch based search application.

http://static.flickr.com/43/117204451_d634c1d869.jpg

3 Implementing a Nutch Search

Implementing our own version of Nutch is fairly easy, provided that you;

1. have a basic understanding of how a web search engine works and
2. are comfortable working in a command line and finally
3. have a fair knowledge of Java and Servlet containers

If you said 'yes' to all three questions above, you have a very high probability of having your Nutch implementation up and running by the end of the steps which follows.

3.1 Before We Begin

3.1.1 Download Nutch

Go to <http://www.apache.org/dyn/closer.cgi/lucene/nutch/> and select a mirror to download Nutch. The version described in this document is version 0.7. After downloading the archive, extract it to your disk.

⚠ NOTE: This document assumes that the archive was extracted to /home/tyrell/nutch-0.7 change this path to reflect your location.

3.1.2 Download and Install a Servlet Container

Apache Tomcat is a popular Open Source servlet container. We will use this to deploy the Nutch search web application in this document. The version referred to in this document is version 5.5. You can download Tomcat from <http://tomcat.apache.org/download-55.cgi>

3.1.3 To Cygwin or not to Cygwin

To run the Nutch shell scripts and create the indexes, we require a UNIX like environment. If you do not have access to a UNIX environment, you can use Cygwin as an alternative. More details and download information for Cygwin can be found at <http://www.cygwin.com/>

3.2 Creating the Index

In order for the nutch web application to function, it will require at least one search index. A search index in nutch is represented in the file system as a directory. However, it is much more than that and is similar in functionality to a database. The nutch API will interact with this index making the internal mechanisms transparent to both developers and end-users.

The steps to create and integrate a new index are as follows;

NOTE: The steps below are assumed to be carried out from inside the /home/tyrell/nutch-0.7 directory created when extracting the archive. Change the path according to your local instance.

3.2.1 Create a directory of root urls.

The nutch 'crawl' command expects to be given a directory containing files that list all the root level urls to be crawled. So create a 'urls' directory in the nutch directory. Then, to crawl the <http://www.virtusa.com> site from scratch, you might start with a file named 'virtusa' in the 'urls' directory, and in the file, add just the URL for the Virtusa home page, <http://www.virtusa.com>. All other pages should be reachable by links from this page.

The 'depth' option to the crawl command will limit how far the crawl goes. Also, the conf/crawl-urlfilter.txt file, described next, will limit what sites to crawl to.

3.2.2 Edit the file conf/crawl-urlfilter.txt

If you are using TRUNK then there is no file called conf/crawl-urlfilter.txt but conf/crawl-urlfilter.txt.template. Just do

```
cat conf/crawl-urlfilter.txt.template|sed 's/MY.DOMAIN.NAME/criaturitas.org/'> conf/crawl-urlfilter.txt
```

If you already have this file then replace the existing domain name with the name of the domain you wish to crawl. For example, if you wished to limit the crawl to the virtusa.com domain, the line should read:

```
+^http://([a-z0-9]*\.)*virtusa.com/
```

This will include any url in the domain virtusa.com in the crawl.

3.2.3 Running a Crawl

Once things are configured, running the crawl is done by using the crawl command.

Its options include:

```
* -dir dir names the directory to put the crawl in.
* -depth depth indicates the link depth from the root page that should be crawled.
* -delay delay determines the number of seconds between accesses to each host.
* -threads threads determines the number of threads that will fetch in parallel.
```

For example, a typical command might be:

```
bin/nutch crawl urls -dir crawl.virtusa -depth 10
```

3.2.4 Output of the crawl

Assuming that the above command is executed with the given parameters, the result will be as follows;

- A new directory will be created named 'crawl.virtusa' in the working directory (according to this guide, /home/tyrell/nutch-0.7/crawl.virtusa)
- This new directory will contain the search index for the URLs given in the flat file named 'urls' (created in the working directory according to this example)
- The 'depth' of the search index will be 10

3.2.5 Errors and Failures

- JAVA_HOME environment variable not set
Set the variable to point to your JDK installation
- Missing 'urls' flat file
Create the file and give the correct path to it in the crawl command
- The indexing domain is not properly defined in the 'crawl-urlfilter.txt', as described in the guide

3.3 Configuring the Nutch Web Application

The search web application is included in your downloaded Nutch archive. In order for the nutch search web application to function properly, it needs to know where to find the indexes. We need to map our indexes by editing the 'nutch-site.xml' file.

The steps to follow would be;

1. Deploy the Nutch web application as the ROOT context
It is not clear why the developers designed the application to run in the root context. However it is possible to modify the application to enable it to be deployed normally.
2. In the web application deployment directory, open the '\WEB-INF\classes\nutch-site.xml' file in a text editor.
3. Change the values of the tags as follows and save the changes.

```
<property>
  <name>searcher.dir </name>
  <value>/home/tyrell/nutch-0.7/crawl.virtusa </value>
</property>
```

4. Re-start Tomcat

3.4 Running a Test Search

Now that we have created the indexes and configured the Nutch web application, the only thing left is to give it a test run.

Open a browser and type your Tomcat URL (ex: <http://localhost:8080>). The following page will greet you if the web application is configured properly.

http://static.flickr.com/42/117204449_d3fe6f8400.jpg

Now type a keyword (ex: virtusa) and click search. If the implementation works as expected, the following results page will be displayed.

http://static.flickr.com/56/117204450_8317279e3a.jpg

3.5 Maintaining Our Index


Now that all is working, we need to think the long term maintenance of the Index. This is a required activity because the web gets updated frequently. New content will appear on sites while existing content might get modified or deleted altogether.

Nutch provides the administrator with a set of commands to update a given index, however performing them manually will not only be tiresome but also unproductive. Since this task need to be carried out periodically it should ideally be scheduled.

3.5.1 Creating a Maintenance Shell Script

Create a new shell script with the commands given in fig 3.4. The script contains three command segments.

1. Commands, which set the environment (ex:JAVA_HOME).
2. Commands to do the index updating.
3. Commands to optimize the updated index.

 **NOTE:** Index optimization is necessary to prevent the index from becoming too large, which will eventually result in a 'too many open files' exception in Lucene.

```
#!/bin/bash

# Set JAVA_HOME to reflect your systems java configuration
export JAVA_HOME=/usr/lib/j2sdk1.5-sun

# Start index updation
bin/nutch generate crawl.virtusa/db crawl.virtusa/segments -topN 1000
s=`ls -d crawl.virtusa/segments/2* | tail -1`
echo Segment is $s
bin/nutch fetch $s
bin/nutch updatedb crawl.virtusa/db $s
bin/nutch analyze crawl.virtusa/db 5
bin/nutch index $s
bin/nutch dedup crawl.virtusa /segments crawl.virtusa/tmpfile

# Merge segments to prevent too many open files exception in Lucene
bin/nutch mergesegs -dir crawl.virtusa/segments -i -ds
s=`ls -d crawl.virtusa/segments/2* | tail -1`
echo Merged Segment is $s

rm -rf crawl.virtusa/index
```

3.5.2 Scheduling Index Updates

The above shell script can be scheduled to be run periodically using a 'cron' job.

