# Requirements for installing Nutch

1. Java 1.4.x, either from Sun or IBM on Linux is preferred. Set NUTCH_JAVA_HOME to the root of your JVM installation.
2. Apache's Tomcat 4.x.
3. On Win32, cygwin, for shell support.
4. Edit the file conf/crawl-urlfilter.txt and replace MY.DOMAIN.NAME with the name of the domain you wish to crawl. For example, if you wished to limit the crawl to the apache.org domain, the line should read:  (This will include any url in the domain apache.org.)

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

(Reference: "Nutch Tutorial", http://lucene.apache.org/nutch/tutorial.html)

Step 1: Perform crawl
```
./nutch crawl ../urls -dir ../crawled/ -depth 1
```
where "urls" file contains one url (http://lucene.apache.org/nutch/) for demo purpose and "crawled" directory is the directory where crawled content will be stored.
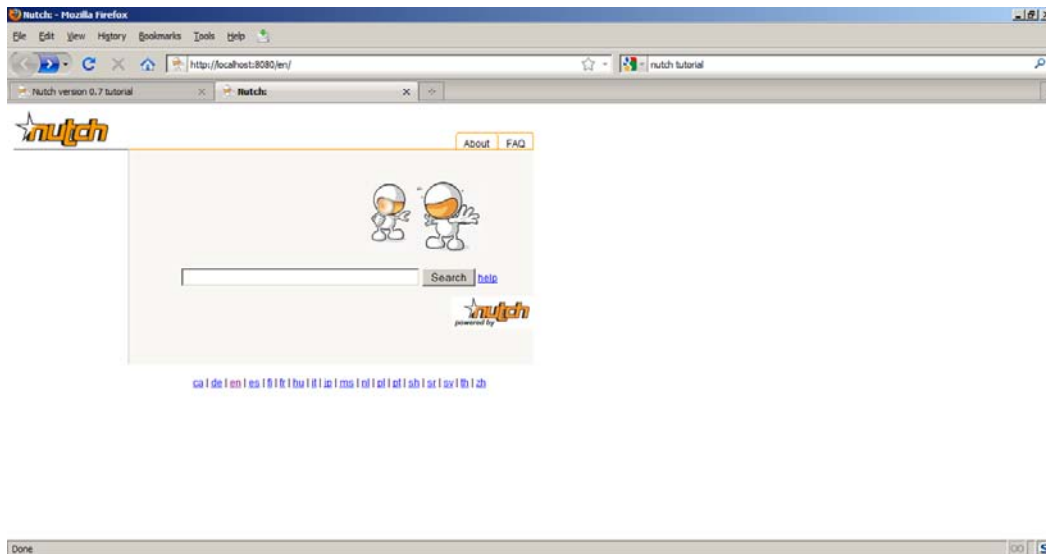
Step 2: Start tomcat server
```
catalina.sh start
```

Step 3: Now open the following URL in a browser to access Nutch search interface
```
http://localhost:8080/
```

The snapshot of this interface is as shown below:

The below snapshot shows the query results for the keyword "apache":