# Assignment 2 : CS 751

HARISH

March 6, 2015

## 0.1 Part 1

**Choose 100 URIs from A1 Generate WARC files of those URIs using:wget,WARCreate,Heritrix (stand-alone or via WAIL),webrecorder.io**

**WARC using webrecorder.io**  I have chosen 100 URIs from Assignment1. In order to generate WARC files in webrecorder.io , we use "https://webrecorder.io/" and record the data using the options provided in the website . The WARC for all the 100 URIs are generated . It can be downloaded in zip format . The file 100interestingurlss.txt contains the 100 urls.

**WARC using WARCreate**  The WARC files using WARCreate can be achieved via chrome plugin 'WARCreate' . The WARC files are downloaded in zipped format (.gz).

**WARC using wget**   I have written a shell script using wget command to download warc files for all the URIs .

**WARC using WAIL**  I have used following steps to set up WAIL in the system.
    1. Install application from webpage, download: Windows 7+, version 0.2013.2.19
2. Unzip
3. Install wx python module: http://www.wxpython.org/download.php (I installed the last one 64-bit, python 2.7)
4. Moved "WAIL..." to C: and rename "WAIL" (remove extra characters)
5. Moved every file folder to the path C:WAIL (That is make direct path to folder, no intermediate folders)
6. Modified BDBCollection.xml as mentioned in google groups.
7. In WAIL.py changed catalina-start.bat .

    I followed the below step for adding multiple URIs to WAIL , to crawl all links at once

    WAIL : Advanced : Setup One-Off Crawl : Click first line to add URI0 : Click first line to add URI1 ... : Clicked Launch Crawl
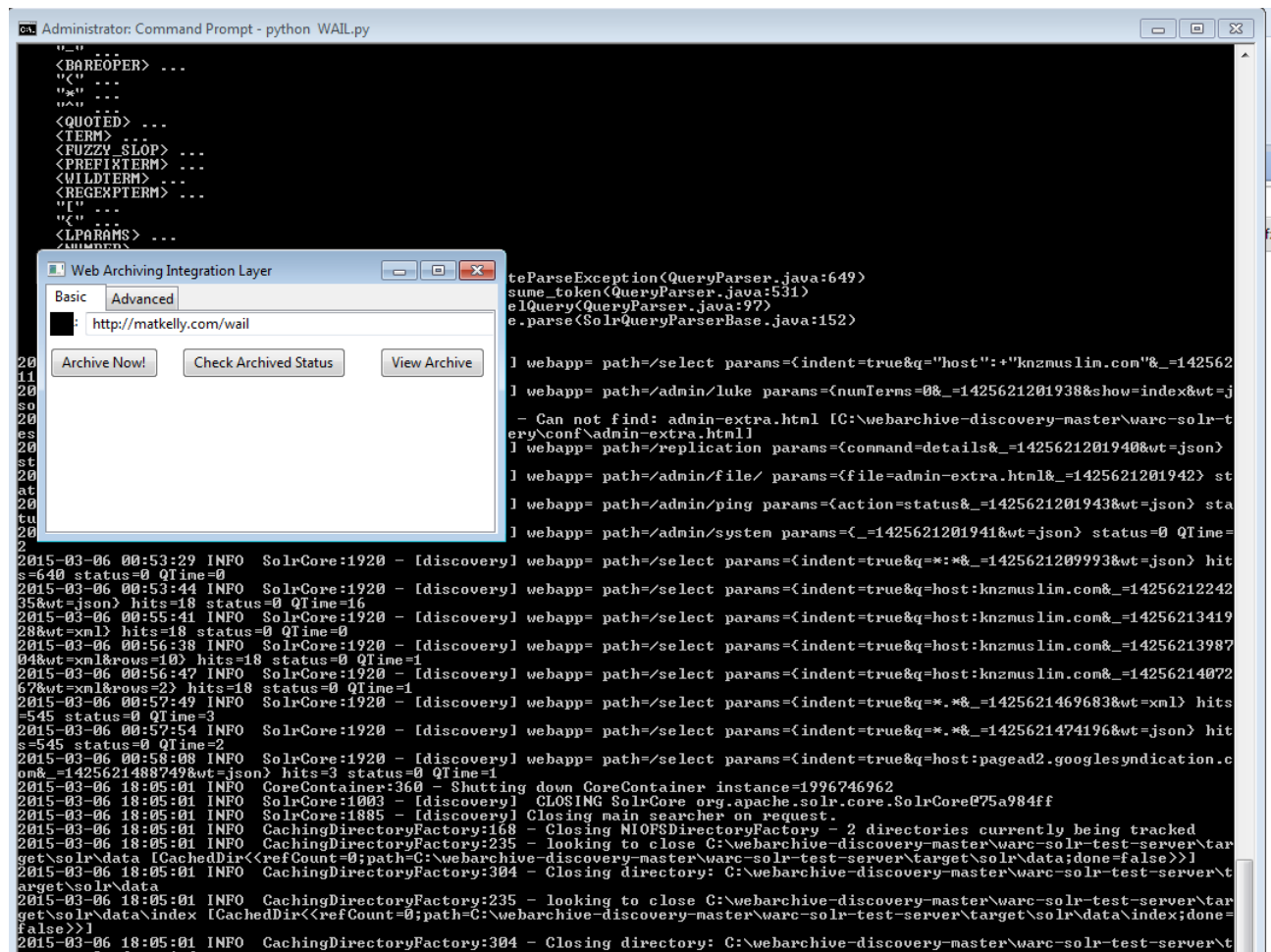
1

## WAIL Instance



Figure 1: Web Archiving Integration Layer

**Describe the resulting WARC files: quantitatively compare and contrast the results of the WARC files of the same URI as generated by different tools** In order to quantitatively compare and contrast the results of the WARC files , I have taken 2 URIs and compared their warc file size generated by 4 different tools.

I have created a bar plot which contains the size differences of the warc files for same URIs. I have also pasted the Rscript(sizeComparision.R located in the same folder) which i have used to plot this graph .

**File Size Comparison**



Figure 2: warc size comparison 1

Figure 3: Histogram 2

The size of same URI warc file depends on the tool we use to generate it . WAIL warc size is comparatively long when compared to other warc file tools because wail crawls for the links present in the URI and gets the data of those links as well. wget warc file size is comparatively smaller than other warc files size.

**Demonstrate playback of 2-3 WARCs in the (Wayback Machine (via WAIL or stand-alone) or pywb) and (webrecorder.io)**  I have played the same WARC file using pywb and webrecorder.io.

I have used this link to setup pywb "https://github.com/ikreymer/pywb ". Pywb takes the warc and uses the following command to generate cdx file.

"cdx-indexer –sort myarchive/cdx myarchive/warcs"

After setting up pywb . Run the wayback server using wayback command .We can search for the URIs using "http://localhost:8080/pywb"

I have included playback screen shots generated using pwyb(Figure 4,5 ) .Also ,I have included the playback screen shots generated using webrecorder.io ( Figure 6 , 7 ,8)

## 0.2 Part 2

**Ingest the 100 URIs from their resulting WARC files into a SOLR instance see the code + tutorial at: https://github.com/ukwa/webarchive-discovery Demonstrate several functioning queries on the files (a full front-end is not required) describe the configuration choices you made in setting up SOLR and processing the documents**

I have set up SOLR instance using the link "https://github.com/ukwa/webarchive-discovery "

The pre-requisites for this are Maven 3 and Oracle java 7 .

Following are the commands used to run SOLR instance.
$cdwarc - solr - test - server$
mvn jetty:run-exploded

This will fire up a suitable Solr instance, with a UI on port 8080 .
Indexing a WARC file

Downloaded warc-indexer from "https://oss.sonatype.org/content/repositories/snapshots/uk/b indexer/2.0.1-SNAPSHOT/"

Copied the downloaded folder to webdiscovery.

Merged all the 100 WARC files using the WARC merger posted in google groups.

Indexed the warc file using the command

–java -jar warc-indexer-2.0.1-SNAPSHOT-jar-with-dependencies.jar -s http://localhost:8080/dis -t warc-indexer/src/test/resources/wikipedia-mona-lisa/flashfrozen-jwat-recompressed.warc.gz

**Querying SOLR** The queries i have used are

1 To search the host name using xml content , change host name data changes automatically (shown in the figure 9 ,17

2 To get json data , set up wt to json (shown in figure 10 )

3 To get xml data , set up wt to xml (shown in figure 16 )

4 To retrieve only two rows from the query , set up row count to 2 (shown in figure 16)

5 To search based on host name in json script(shown in figure 10 )

This is

Products
Customer Support
Shipping/Returns/Exchange
Contact
Cart



PROCEE

Products
Cart
More

100 Emojis Joggers
$35.00

localhost:8080/pywb/20150304182926/https://newteachersblog.wordpress.com/2010/11/11/professional-bounda

Apps    Customize Links

○ gtceblog
January 11, 2011 at 6:00 pm

Thanks for that Shauna. So, in your view, is there no limit in what a teacher can do to him or herself? and still retain y
Just trying to see where all our boundary limits are – if at all. Interesting stuff. Thanks.

Reply

3. Aaron Puley
January 31, 2011 at 4:47 am

This is definitely a topic I feel passionate about. I am a teacher, and work diligently to kept at the forefront of teaching pract
tattoos up both arms in complete sleeves. I have a back tattoo in the planning and the images of my children are in considera
conservative in nature (although not believing it myself), I removed 4 earrings from each ear to establish the "acceptable" in
So, I came up with a happy medium to remove half of my earrings per ear for the interview. That being said, I also cut my h
change yet knew it was necessary to make the hire. Being comfortable with oneself also means understanding the expectatio
system with a "us versus them" mentality, but as an adult, a deeper understanding of the system and how to work within it i

Flash forward 6 years, and I can tell you in retrospect that I encountered a progressive principal and he was more interested
person by the length of their hair or the metal in their ears or the ink on their arms is no different than judging them for the c
the scope of their child's education as result? Definitely. I have 2 undergraduate degrees, a Master's degree, and a Bachelor
innovative curriculum top students every month of every year.

In the early years of my teaching practice, I wore long sleeved shirts to work. Students saw the peak of my tattoos while I w
them, however, that it is important to understand that some people do not understand or appreciate the differences of others

Individuality has always come at a price.

I can honestly say, that no parent has looked at me as if I am deficient in the teaching of their child. Parents want to know th

Reply

○ gtceblog
January 31, 2011 at 10:16 am

Aaron, thanks for that excellent post. You cover a lot of ground there which I am sure will be food for thought for a l
other people's sensitivities (and even prejudices) – which in my view is an excellent thing in a teacher. I am sure you

But here's one more scenario for you to consider, or rather two. First, let's say a female teacher colleague for exampl
or principal of the school. Do you say anything?

Secondly, your child's new teacher is a devout Muslim who chooses to cover her body from head to toe including her

You don't have to post an answer to these if you don't want to – as you can see – I'm just pushing the 'hypotheticals'

Reply

4. Pingback: Can you be a Good Teacher with Tattoos and Body Piercings? | The Creative Education Blog

5. Pingback: Tweets that mention Professional boundaries… where does it all end? « gtceblog -- Topsy.com

6. Pingback: Professional boundaries… where does it all end? « gtceblog | Max's World

10

7. Gadget
June 13, 2011 at 8:09 pm

I'm a secondary teacher in my 3rd year of teaching. When I went into teaching I had one tattoo now upto 5 of which in the s
I'm also in the current process of growing dreadlocks and radically changed my appearance by cutting all my hair off to beg
This has been well received by the kids I teach, and haven't had any of the assumptions or rude comments I was possibly ex
conform.
The scenario of dress code, in my opinion does come down to professionalism I wear, bright colours skirts and dress not dre

Figure 6: webrecorder

Figure 7: webrecorder plaback file

Figure 8: webrecorder playback file 2

Figure 9: Querying solrInstance

Figure 10: Querying solrInstance

15

Figure 11: Querying solrInstance

Figure 12: Querying solrInstance