# CS 751: Assignment 2

Bharath Kongara

Spring 2015

# Contents

# 1 Question 1

Choose 100 URIs from Assignment1 and generate WARC files of those URIs using:

- wget

- WARCreate

- Heritrix (stand-alone or via WAIL)

- and webrecorder.io

Describe the resulting WARC files: quantitatively compare and contrast the results of the WARC files of the same URI as generated by different tools

- choose interesting examples

Demonstrate playback of 2-3 WARCs in the (Wayback Machine (via WAIL or stand-alone) or pywb) and (webrecorder.io)

- "`https://github.com/iipc/openwayback`"

- "`https://github.com/ikreymer/pywb`"

## 1.1 Solution

The following steps were taken to setup tools and generate WARC files:

- Installed wget using 'brew install wget'.

- WARC file for each URI is generated using the command 'wget –warc-file=outputfilename URI'.

- Downloaded WARCreate chrome extension from chrome web store and added it to the chrome browser. Generated WARC for each URI manually by providing input to WARCreate.

- Installed WAIL.

- The figure below shows how to generate WARC files for multiple URIs using single heritrix instance.
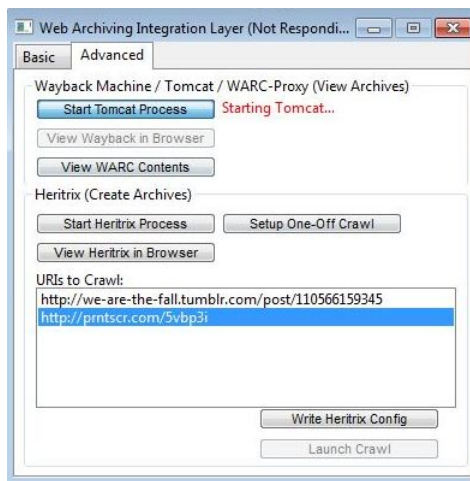


Figure 1: Generating WARC files for multiple URIs using WAIL

- Of the 100 URI selected, when any of the URIs had a parameter then WAIL wasn't able to generate WARC file for it.

- Used "`https://webrecorder.io/`" . Generating WARC files for multiple URIs is done as shown in the figure below.
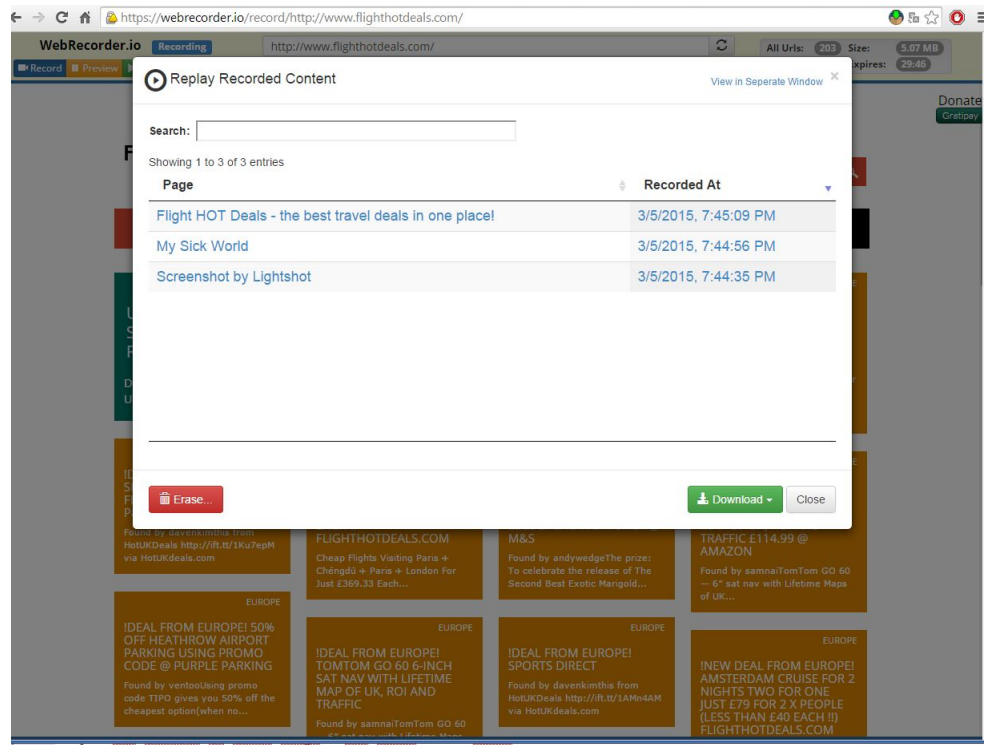
Figure 2: Generating WARC files for multiple URIs using webrecorder

- Quantitative comparison of WARC files generated by WAIL,webrecorder,WARCreater and wget is done on the basis of WARC file sizes.

- The comparison sizes are WAIL = 40 MB , webrecorder.io = 20.48 MB, WARCreate = 15.36 MB, wget = 4 MB.

- The WAIL size is more compared to others because WAIL software crawls data of the links in the website.

- This comparison is shown in the graph below.

- Installed pywb to playback WARC files.

- pywb requires cdx for each WARC to playback, so generated .cdx file for each WARC file.

- I have put all my WARC files into a new folder and changed the archive paths of config.xml to point to my WARC files.

- Play back for two WARC files using pywb is shown in the below figure.

- The WARC file which was in the play back does not have some images and some URLs.

- The WARC file which is in the play back does not contain certain CSS archives and some links

- WebRecorder.io takes WARC file to replay the archived file. Play back for two WARC files using WebRecorder.io is shown in the below figure

- The WARC file in this play back worked perfectly without any faults.

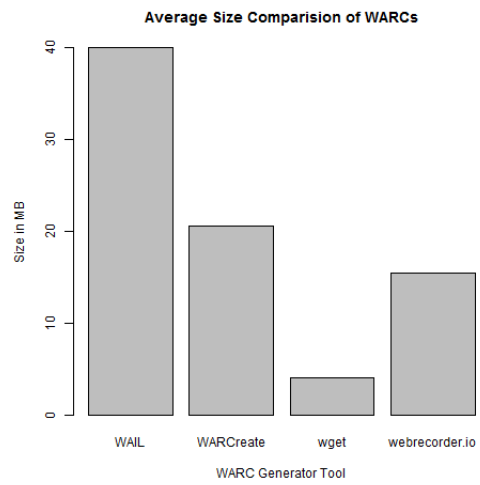- The WARC file in this play back did not have certain images.

Figure 3: Average Size Comparision of WARC
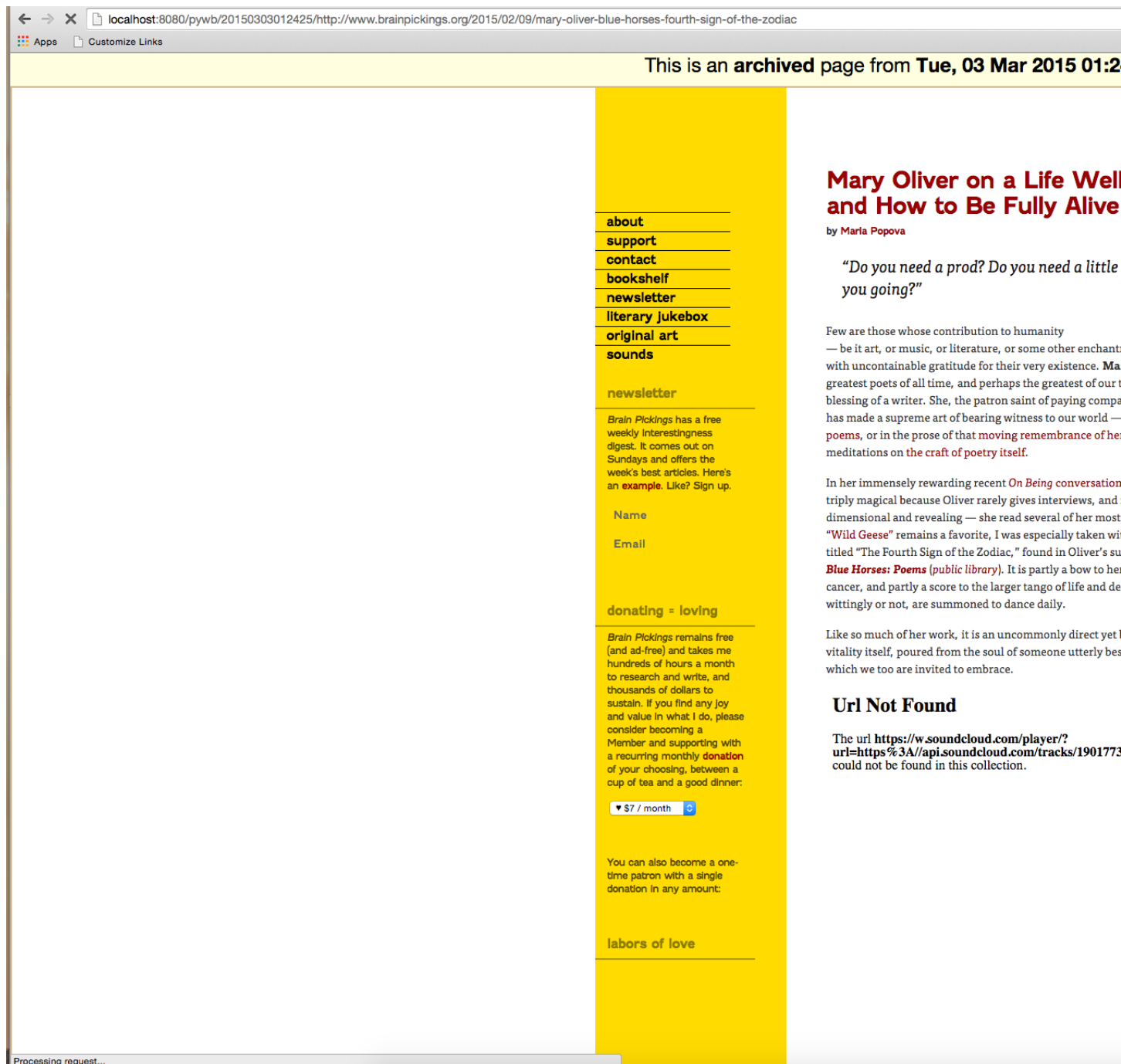
Figure 4: Play back WARC file using pywb

Figure 5: Play back WARC file using pywb
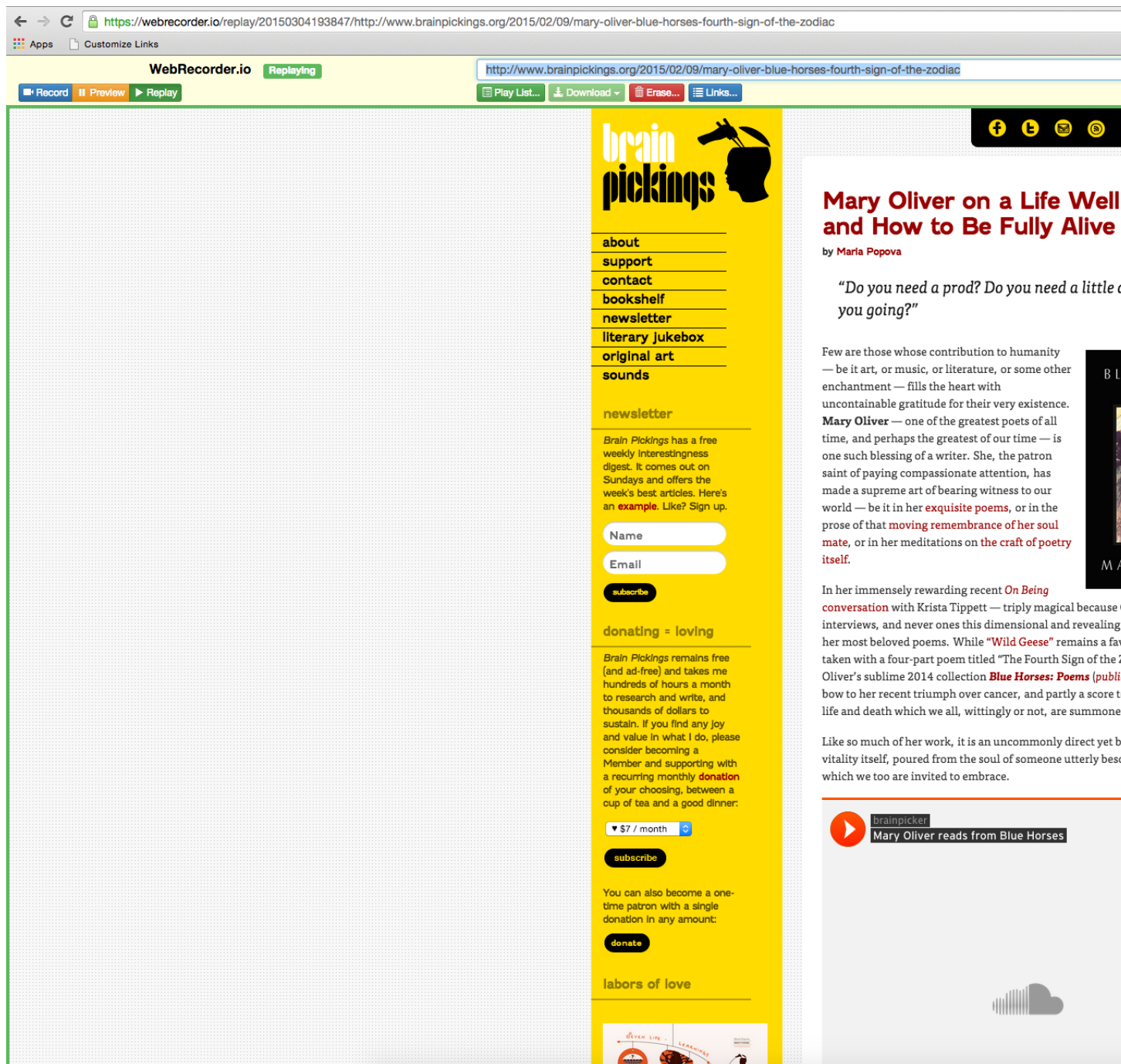
Figure 6: Play back WARC file using webrecorder

Figure 7: Play back WARC file using webrecorder

# 2 Question 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

## 2.1 Solution

The following steps were taken to configure SOLR and process documents :

- The pre-requisites for SOLR are Maven 3, Java 7 so I installed them.

- I faced an issue while installing SOLR which is jetty dependency not found.

- I got it working by adding the dependency to the pom.xml

- The command "mvn jetty:run-exploded" starts the SOLR instance.

- I indexed the WARC file by using the command "java -jar path of jar -s "`http://localhost:8080/discovery`" -t Path of WARC".

- I tried to set-up Shine for SOLR as the front end instead of the default UI for SOLR but after trying for two days and multiple email exchanges with the author of the tool I wasn't able to. So I used the default SOLR front end.

- In SOLR we can perform the following queries.

  - Here we are demonstrating how to retrieve the names and ids of all documents with "`http://localhost:8080/solr/select?q=inStock:false&wt=json&fl=id,name`"

  - Here we are using the functional query idf(field,term). This function returns the inverse document frequency for the given term, using the similarity for the field. "`http://localhost:8080/solr/select/?fl=score,id&defType=func&q=mul(tf(text,memory),idf(text,memory))`"

  - The functional query being used here is tf(field,term) and it returns the inverse document frequency factor for the given term using the similarity for the field. "`http://localhost:8080/solr/select/?fl=score,id&defType=func&q=mul(tf($f,$t),idf($f,$t))&f=text&t=memory`"

  - termfreq(field,term) is the functional query that is being used and it returns the number of times the term term appers for that field in the document. "`http://localhost:8080/solr/select/?fl=score,id&q=DDR&sort=termfreq(text,memory)desc`"

  - Here norm(field) is the functional query that is being used. It returns the norm stored in the index, the product of the index time boost and then length normalization factor. "`http://localhost:8983/solr/select/?fl=score,id&q=DDR&sort=norm(text)asc`".

- Below are the snapshots of the queries that I ran on SOLR after indexing.

Figure 8: Query 1 on SOLR



Figure 9: Query 2 on SOLR

Figure 10: Query 3 on SOLR

# Bibliography

[1] Download warcreate. https://chrome.google.com/webstore/detail/warcreate/kenncghfghgolcbmckhiljgaabnpcaaa?hl=en-US.

[2] Generating simple bar plots in r. http://www.statmethods.net/graphs/bar.html.

[3] Playback using webrecorder. https://webrecorder.io/.

[4] Setting up pywb. https://github.com/ikreymer/pywb.

[5] Using wail. http://matkelly.com/wail/.

[6] Warc merge. https://github.com/maturban/WARCMerge.