

Khaja Masroor Ahmed

Assignment 2

CS 751: Introduction to Digital Libraries

Dr. Michael Nelson

Spring 2015

March 6, 2015

Contents

1	Question 1	1
	1.1 Question	1
	1.2 Solution	1
	1.3 Comparison	2
	1.4 Playback	7
	1.4.1 Webrecorder.io	7
	1.4.2 pywb	19
	1.5 Code Listing	26
2	Question 2	27
	2.1 Question	27
	2.2 Solution	27
	References	33

Question 1

1.1 Question

- Choose 100 URIs from A1.
- Generate WARC files of those URIs using:
 - wget
 - WARCreate
 - Heritrix (stand-alone or via WAIL)
 - webrecorder.io
- Describe the resulting WARC files: quantitatively compare 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>

1.2 Solution

- Fetching the 100 URIs was an arduous task, as a lot of the URIs had explicit content or were not accessible anymore. I had to open most of the links manually and as I progressed through the links I noticed a pattern and started avoiding those links to avoid additional time fetching the URIs.
- I saved the URIs into a uri.txt file.
- Soon after I installed the plug-in for WARCreate from Chrome Web Store and started manually visiting these URIs and clicking on the Warccreate button in the address bar for chrome web browser for downloading the WARC files.
- After WARCreate, I moved onto the next tool webrecorder.io. I manually opened each of the URIs on the website and started downloading the

WARC files. While fetching them manually I realized this activity could be automated but I decided that performing this task manually would re-assure me that the WARC files have been created properly.

- The third tool WAIL consumed most of my time as the set-up wasn't as simple as it had multiple dependencies, required configuration changes and lacked proper documentation. As the thought of setting up stand-alone heritrix crept into my mind, WAIL set-up issues were addressed by an email from Alexander Nwala and after following his instructions I was able to get the WAIL tool working.
- I used the WAIL tool to create the archive but it took too long for most of the links so I put a time limit of approximately 10 minutes for each of the URIs and then manually stopped execution and moved to the next URI.
- "wget" was the simplest of all the tools for WARC creation. I used a simple unix command with the template "wget URI -warc-file=FILENAME -no-warc-compression" for each of the 100 URIs.

1.3 Comparison

- From the 10000 URIs retrieved from the first assignment we took a sample of 100 URIs. To make a meaningful comparison of the sample I calculated the mean of the size of WARC files from each of the tools.

- Below is the graph to represent the average size of the WARC files for each of the tools.

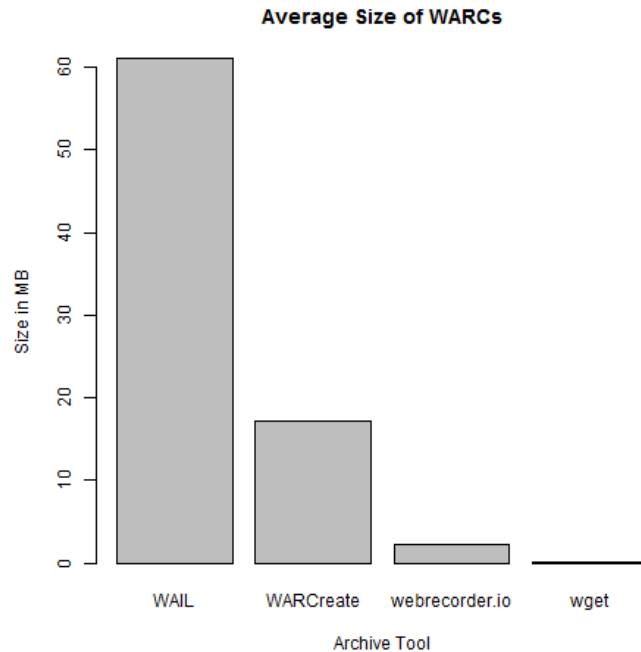


Fig. 1.1. Average Size of WARC file for sample size=100

- While downloading the WARC files for each of the above tools I noticed a trend that WAIL was able to crawl the deepest and it had to be manually stopped. The size of the files generated by WAIL were significantly larger than the ones created by the other tools.
- With this information and the above graph I concluded that WAIL creates the largest of the files. Upon further analysis I found that WAIL crawls dereferences each of the URIs in the webpage. WAIL does not stop until it crawls to the end of the URI.
- As it is a configurable tool, we can modify the specify the time spent to crawl each URI, maximum size of the WARC file and/or maximum documents downloaded in the crawler-beans.xml created for each of the jobs.
- Below are the WARC size plots for three randomly chosen URIs as specified below, which I've used for the rest of the assignment for playback demonstration.
 - URI#1: <http://edition.cnn.com/videos/us/2015/02/04/lead-dnt-johns-plane-crash-taiwan.cnn>
 - URI#2: <http://www.dccomics.com/blog/2015/02/04/breaking-news-milo-ventimiglia-is-coming-to-gotham>

4 1 Question 1

- URI#3: <http://www.theguardian.com/football/2015/feb/04/jonny-evans-manchester-united-louis-van-gaal>

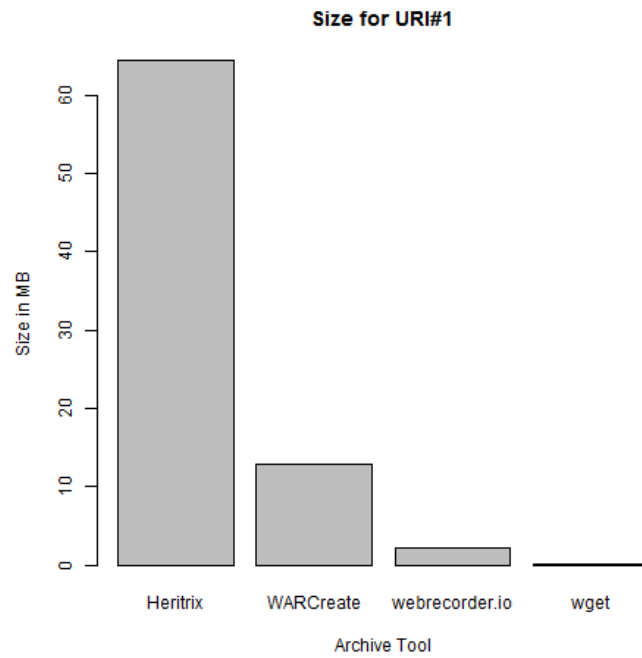


Fig. 1.2. Archive Tools vs. Size for URI#1

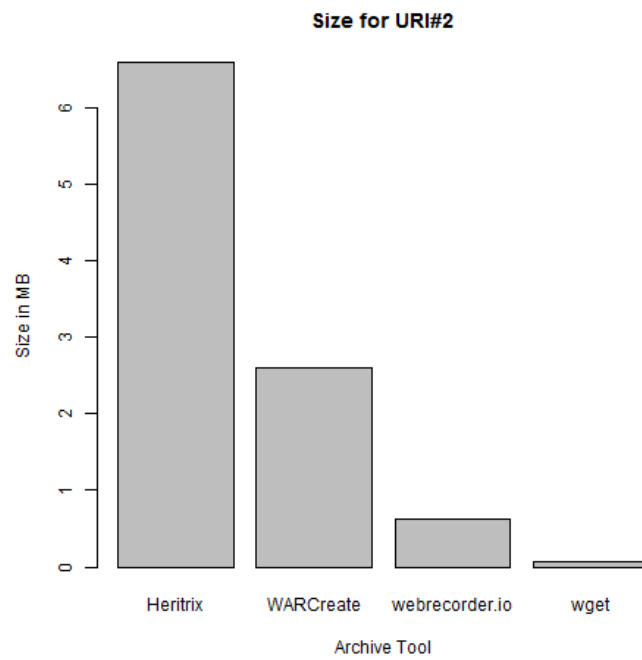


Fig. 1.3. Archive Tools vs. Size for URI#2

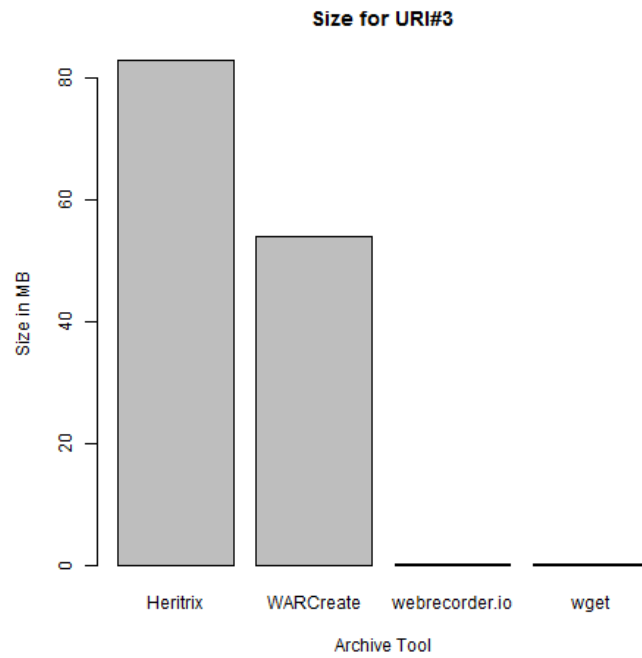


Fig. 1.4. Archive Tools vs. Size for URI#3

1.4 Playback

1.4.1 Webrecorder.io

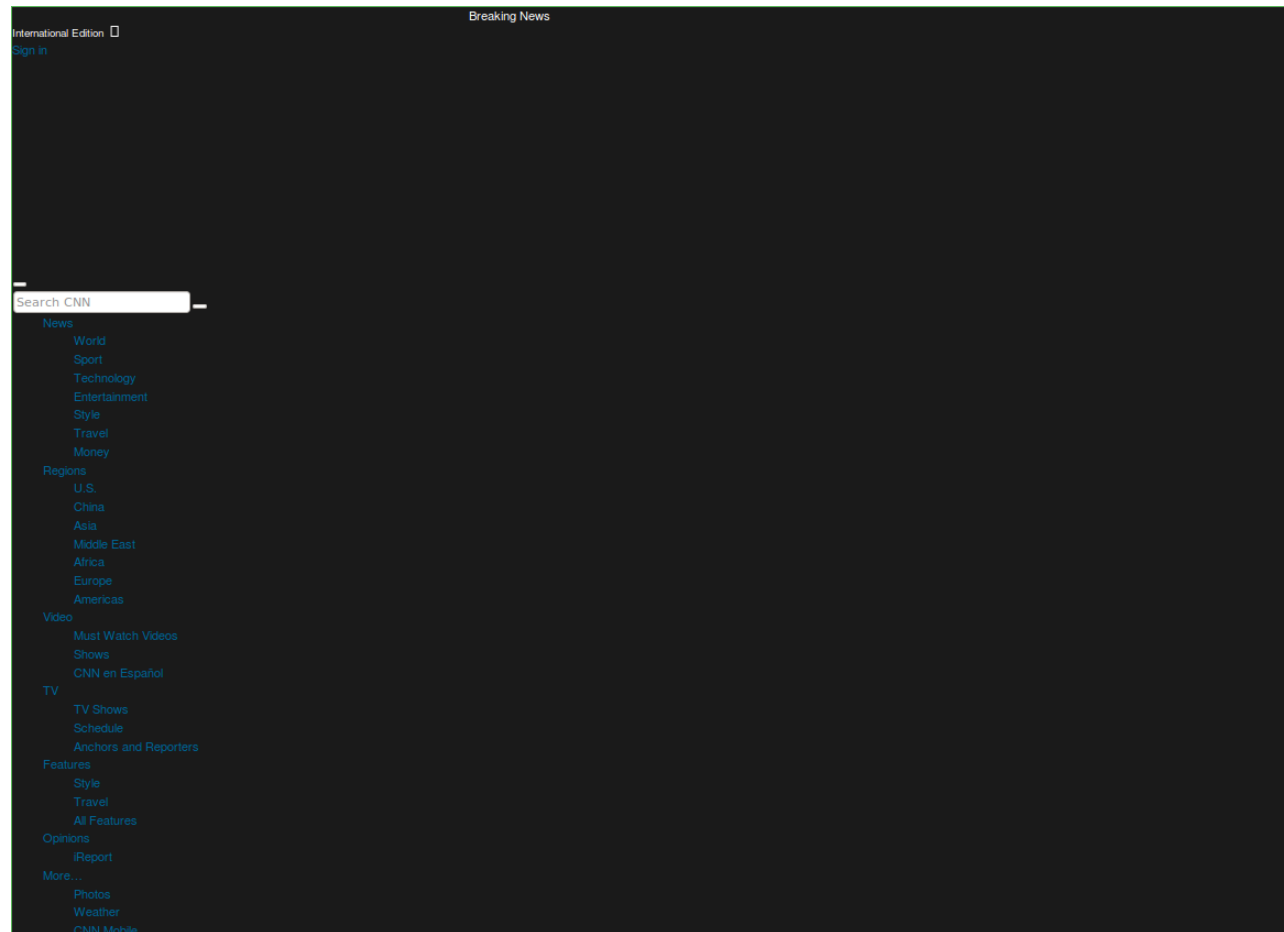


Fig. 1.5. Web Recorder - WAIL generated WARC for URI#1

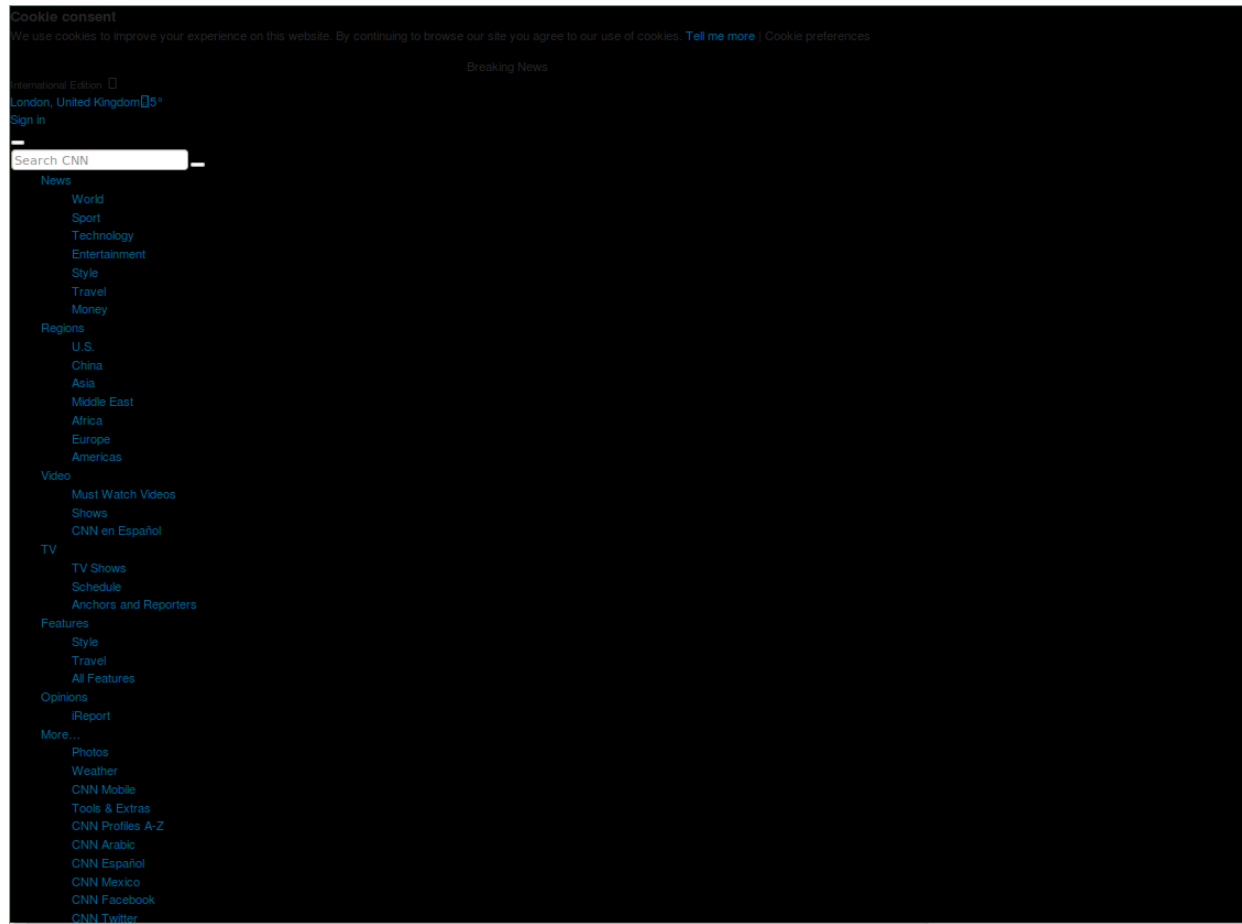


Fig. 1.6. Web Recorder - WARCcreate generated WARC for URI#1

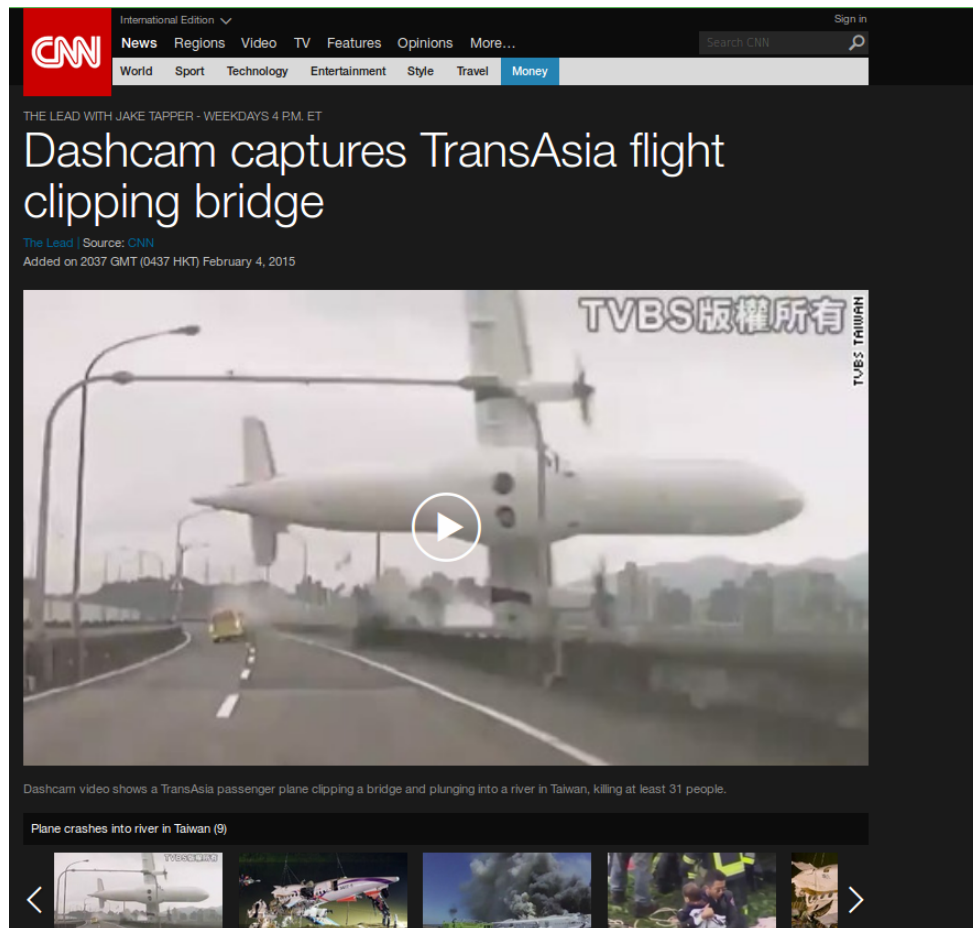


Fig. 1.7. Web Recorder - webrecorder.io generated WARC for URI#1

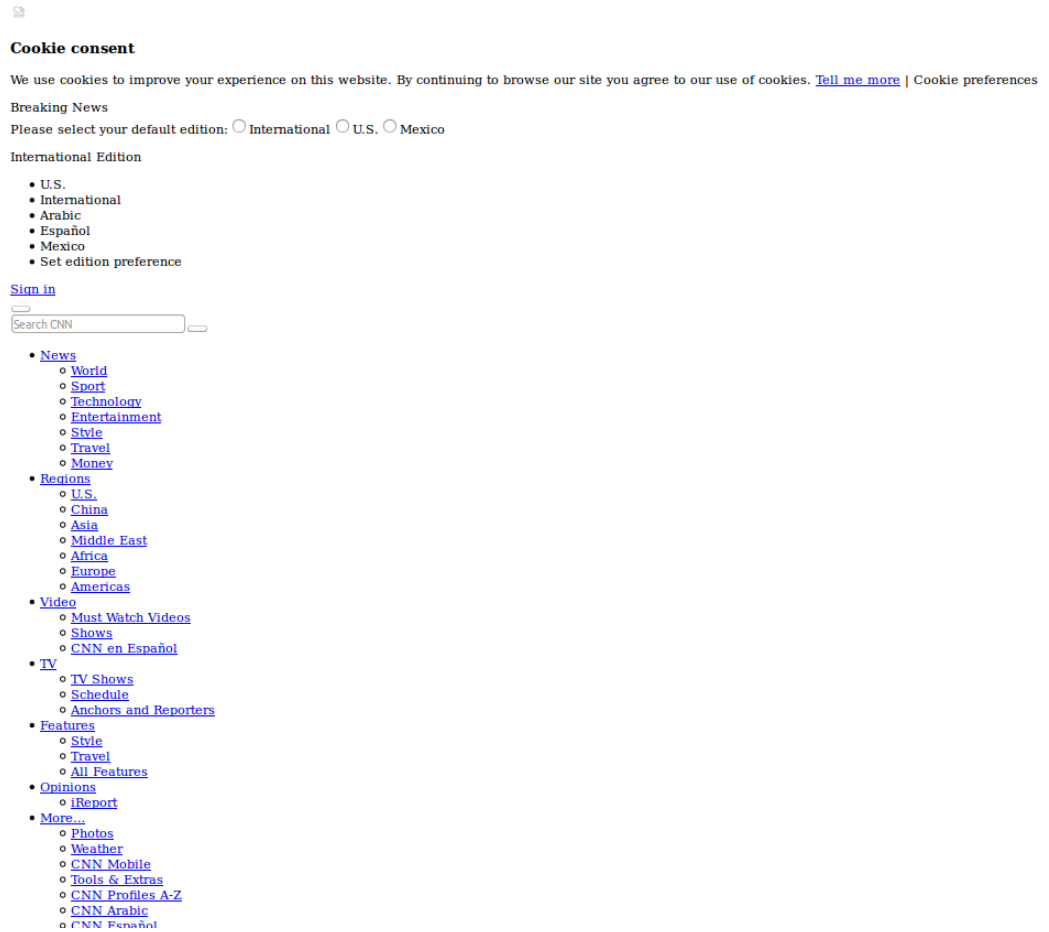


Fig. 1.8. Web Recorder - wget generated WARC for URI#1



He may have played a hero in the past, but this spring, Milo Ventimiglia is getting in touch with his villainous side.



Fig. 1.9. Web Recorder - WAIL generated WARC for URI#1

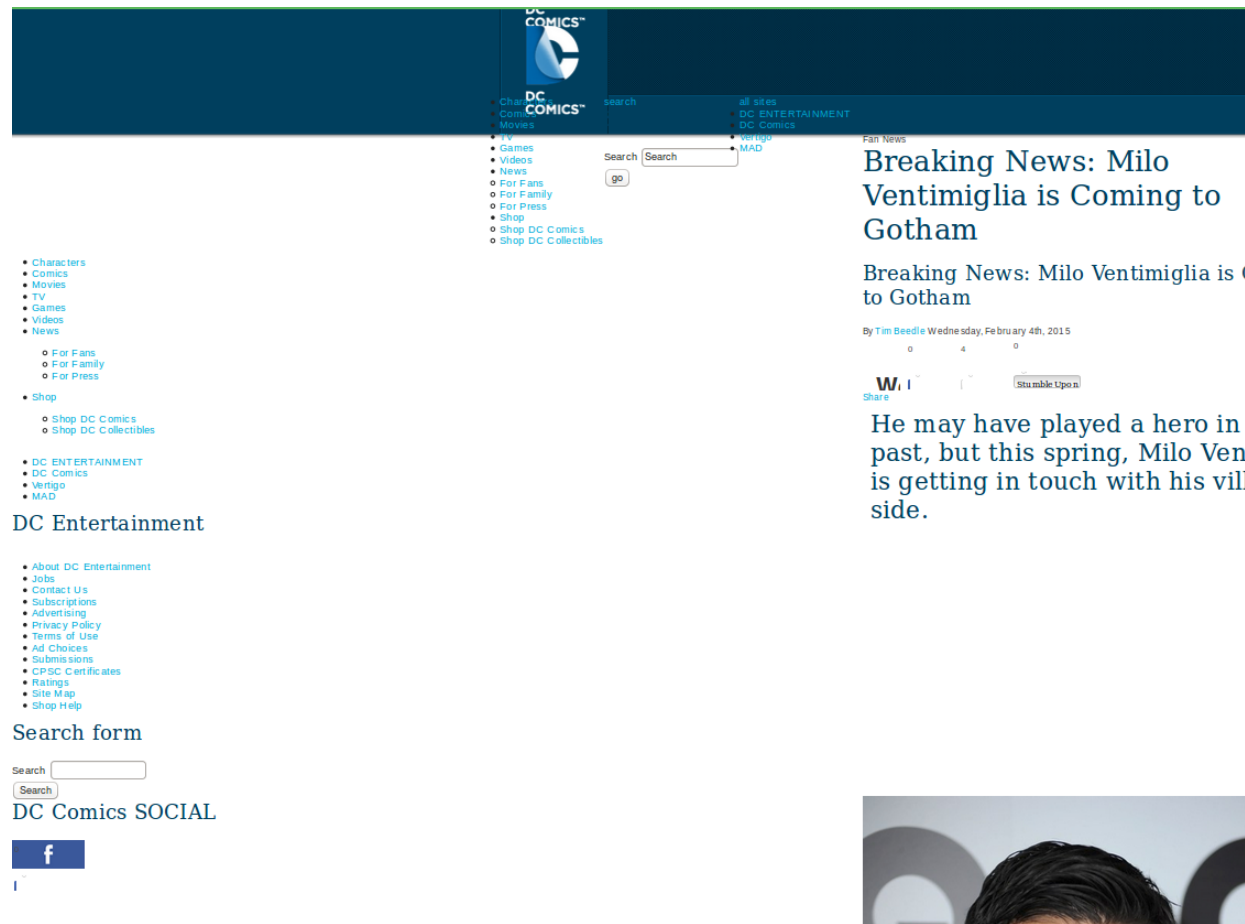


Fig. 1.10. Web Recorder - WARCcreate generated WARC for URI#1



He may have played a hero in the past, but this spring, Milo Ventimiglia is getting in touch with his villainous side.



Fig. 1.11. Web Recorder - webrecorder.io generated WARC for URI#1

- [DC ENTERTAINMENT](#)
- [DC Comics](#)
- [Vertigo](#)
- [MAD](#)

DC Entertainment

- [About DC Entertainment](#)
- [Jobs](#)
- [Contact Us](#)
- [Subscriptions](#)
- [Advertising](#)
- [Privacy Policy](#)
- [Terms of Use](#)
- [Ad Choices](#)
- [Submissions](#)
- [CPSC Certificates](#)
- [Ratings](#)
- [Site Map](#)
- [Shop Help](#)

Search form

Search

DC Comics SOCIAL

Fan News

Breaking News: Milo Ventimiglia is Coming to Gotham

Breaking News: Milo Ventimiglia is Coming to Gotham

By [Tim Beedle](#) Wednesday, February 4th, 2015
[Share](#)

He may have played a hero in the past, but this spring, Milo Ventimiglia is getting in touch with his villainous side.

He may have played a hero in the past, but this spring, Milo Ventimiglia is getting in touch with his villainous side.

[Gotham's](#) already seen its fair share of gangsters, hit men and full on psychopaths, but Detective Jim Gordon will find himself going toe-to-toe with a new kind of serial killer later this season p
multiple-episode recurring arc.

Fig. 1.12. Web Recorder - wget generated WARC for URI#1



Fig. 1.13. Web Recorder - WAIL generated WARC for URI#1

close
next
prev
info-button

The screenshot shows a web browser interface with a dark blue header. The header contains navigation links: 'sign in', 'subscribe', and 'search'. Below the header is a secondary navigation bar with links for 'US', 'world', 'opinion', 'sports', 'soccer', 'tech', 'arts', 'lifestyle', 'fashion', 'business', 'money', 'travel', and 'environment'. The main content area is divided into two columns. The left column features the 'Manchester United' logo and a sidebar for the author, 'Jamie Jackson', including his Twitter handle '@JamieJackson__' and a share count of 540. The right column displays the article title 'Jonny Evans: Manchester United are slowly adapting to Louis van Gaal's style'. Below the title is a list of three bullet points: 'The lads have taken a lot of board. At times it has been tough', 'The FA Cup looks our best chance of a trophy', and 'Match report: Manchester United 3 Cambridge United 0'. A large photograph of Jonny Evans in a red Manchester United jersey is shown. Below the photo, a caption reads: 'Manchester United's Jonny Evans says the players are beginning to adapt to the Louis van Gaal methodology. Photograph: Kieran McManus/BPI/REX'. The article text begins with 'Jonny Evans has admitted Louis van Gaal's new style can be tough for Manchester United to adapt to, though the defender believes it will prove beneficial in the long term.' and continues with 'The Northern Irishman returned to the team after more than a fortnight out with a foot injury for Tuesday's FA Cup victory over Cambridge United. The 3-0 win over the League Two side took United through to face League One's Preston North End in the fifth round and United are third in the'.

Fig. 1.14. Web Recorder - WARCcreate generated WARC for URI#1



Fig. 1.15. Web Recorder - webrecorder.io generated WARC for URI#1

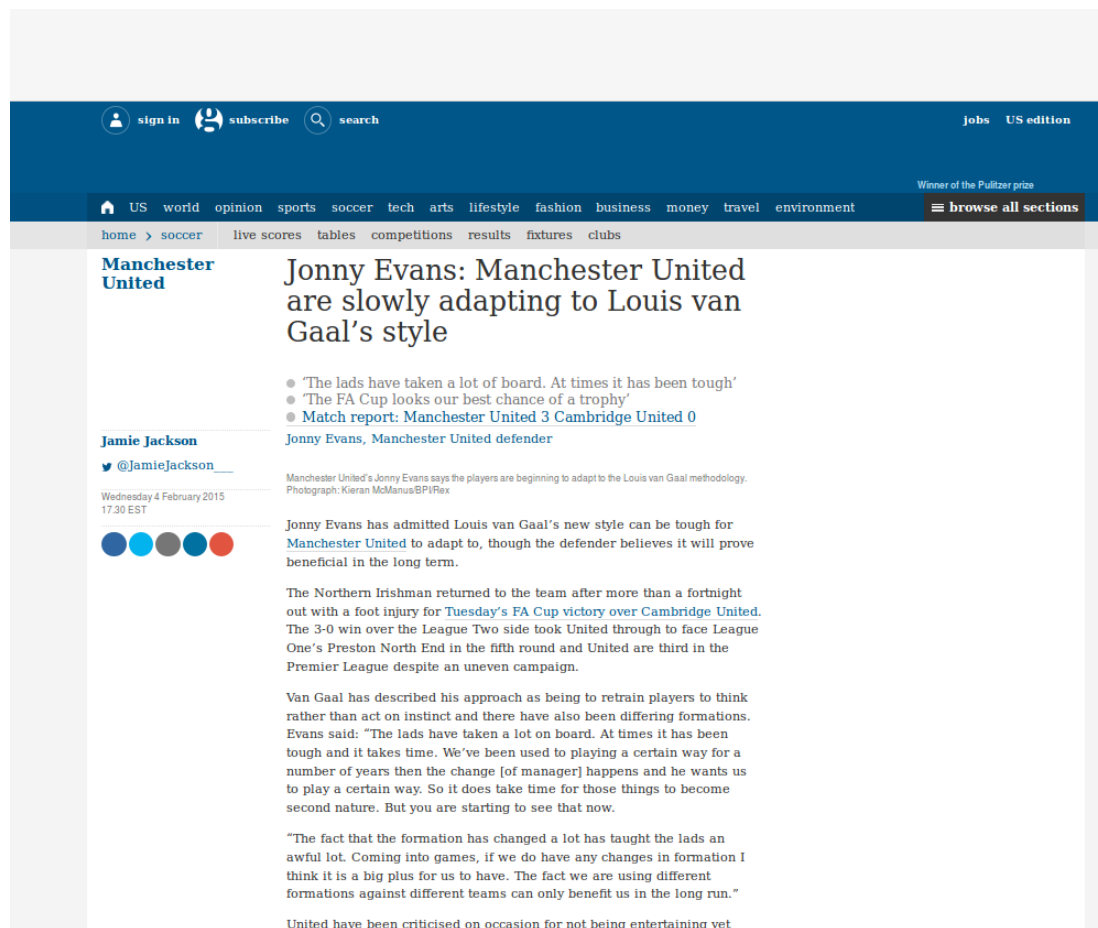


Fig. 1.16. Web Recorder - wget generated WARC for URI#1

1.4.2 pywb

- For setting pywb I followed the instructions at <https://github.com/ikreymer/pywb>.
- I created the cdx files for the WARC files and then launched pywb using the command `wayback` after going to the root folder of pywb in the terminal.
- I used the web browser and went to "localhost:8080" for accessing the web interface for pywb. Below is the search interface for pywb.

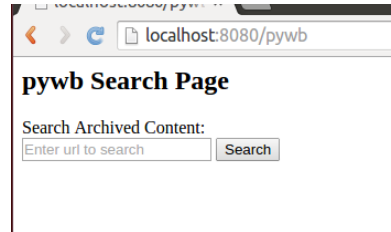


Fig. 1.17. Web Search Interface pywb

- After creating the cdx files for the associated URIs, enter the URI in the search field.
- I was facing an issue with the WAIL results not being displayed for pywb inspite of the cdx files being generated properly and no errors in the command prompt. I tried using different WAIL generated WARC files but I was still not able to view it in the list of search results and hence wasn't able to display the results.



Fig. 1.18. Search results page on pywb

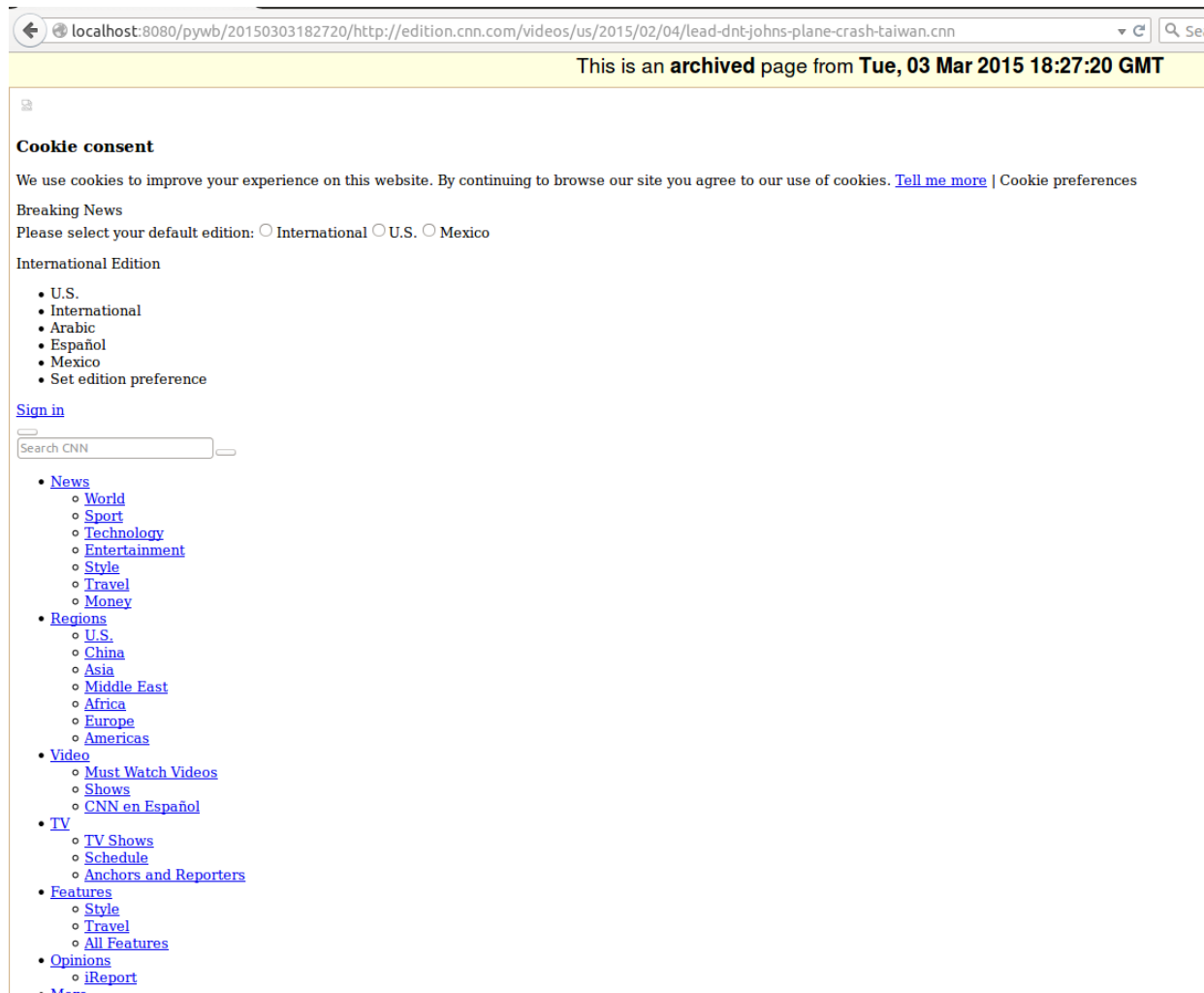


Fig. 1.19. pywb - wget generated WARC for URI#1



Fig. 1.20. pywb - WARCreate generated WARC URI#2



He may have played a hero in the past, but this spring, Milo Ventimiglia is getting in touch with his villainous side.

He may have played a hero in the past, but this spring, Milo Ventimiglia is getting in touch with his villainous side.

Gotham's already seen its fair share of gangsters, hit men and full on psychopaths, but Detective Jim Gordon will find himself going toe-to-toe with a new kind of serial killer later this season played by Ventimiglia (*Heroes*), who will guest star on the hit show in a major, multiple-episode recurring arc.

Ventimiglia will be playing an original character named Jason Lennon, who also goes

Fig. 1.21. pywb - webrecorder.io generated WARC URI#2



He may have played a hero in the past, but this spring, Milo Ventimiglia is getting in touch with his villainous side.

He may have played a hero in the past, but this spring, Milo Ventimiglia is getting in touch with his villainous side.

Gotham's already seen its fair share of gangsters, hit men and full on psychopaths, but Detective Jim Gordon will find himself going toe-to-toe with a new kind of serial killer later this season played by Ventimiglia (*Heroes*), who will guest star on the hit show in a major, multiple-episode recurring arc.

Ventimiglia will be playing an original character named Jason Lennon, who also goes

Fig. 1.22. pywb - wget generated WARC URI#2

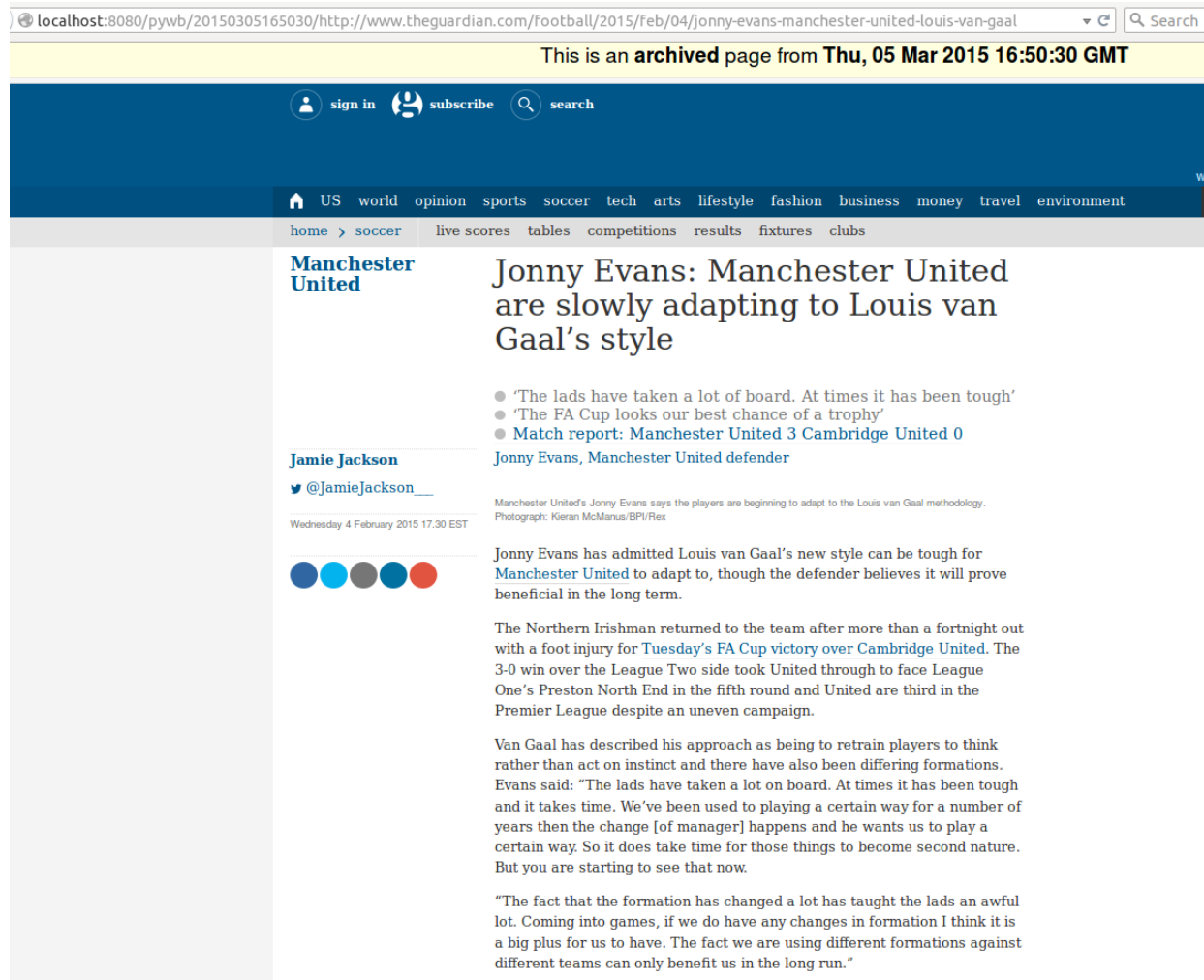


Fig. 1.23. pywb - webrecorder.io generated WARC URI#3

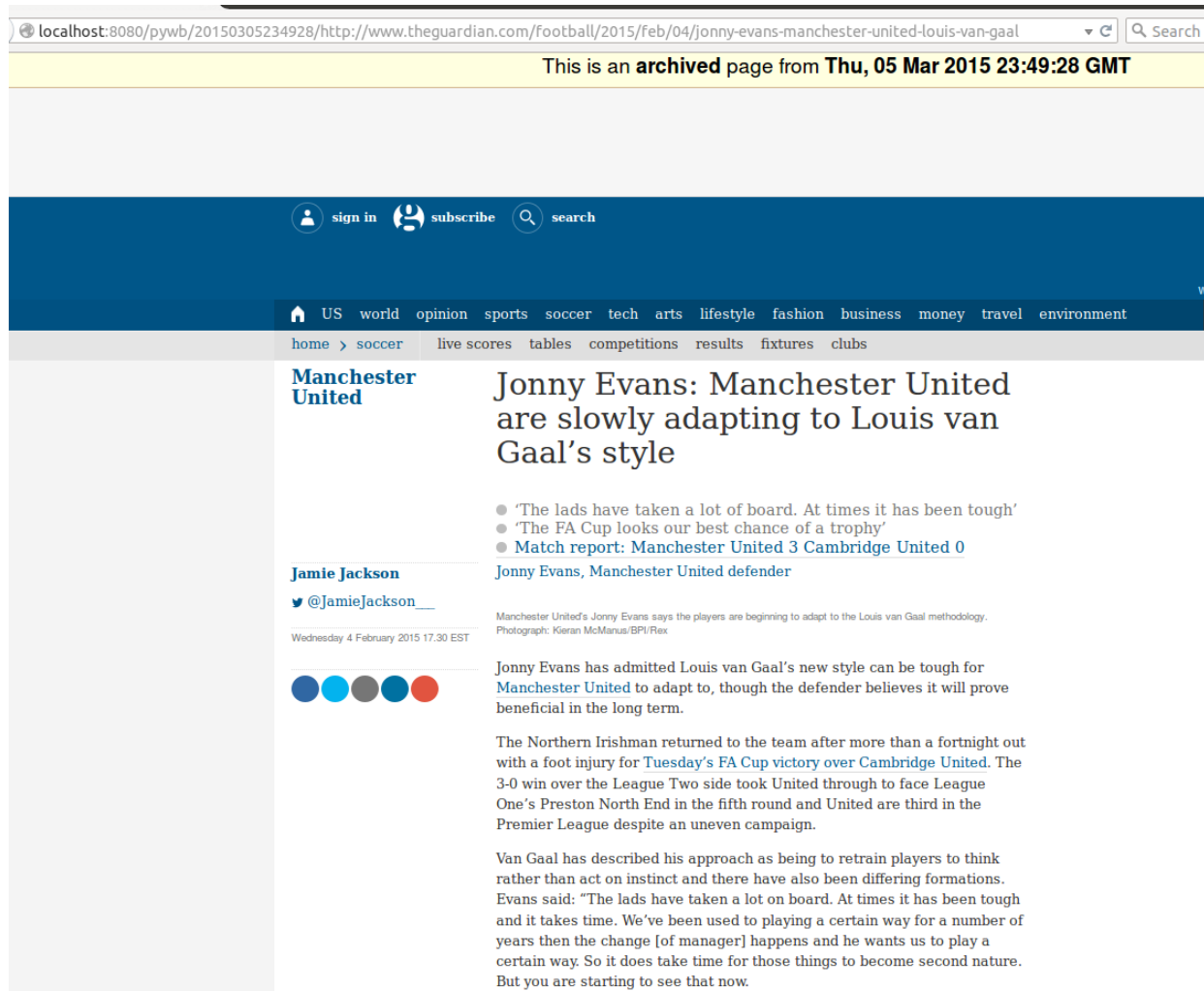


Fig. 1.24. pywb - wget generated WARC URI#3

1.5 Code Listing

```
1 B <- c(83.1,54,0.1537,0.240)
2 png("C:/Users/kahmed/Desktop/uri3.png")
3 barplot(B, main="Size for URI#3", xlab="Archive Tool", ylab=
  "Size in MB", names.arg = c("Heritrix", "WARCreate", "
  webrecorder.io", "wget"))
```

Listing 1.1. R program for generating the bar plot for WARC Tools vs. Average Size of WARC files for 100 samples

Question 2

2.1 Question

- 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

2.2 Solution

- For setting up SOLR I used the link as provided in the question.
- The instructions provided were clear and concise which helped me set-up the necessary tools without any issues.
- Though I faced an issue when the pywb was running and I wouldn't be able to run SOLR. Then I realized that both the applications were trying to run on the same port as it was already in use. Upon closing the pywb instance and then running SOLR everything was working like a well oiled machine.
- I merged all the WARC files that I had retrieved using wget, using the tool WARCMerge located at <https://github.com/maturban/WARCMerge>.
- After merging the WARC file I ran the command for indexing as provided in the user instructions by the author of the tool.
- Following this I ran the default web interface for SOLR and ran the following queries as displayed below.

localhost:8080/#/discovery/query

- Dashboard
- Logging
- Core Admin
- Java Properties
- Thread Dump
- discovery
 - Overview
 - Analysis
 - Dataimport
 - Documents
 - Files
 - Ping
 - Plugins / Stats
 - Query
 - Replication
 - Schema Browser

Request-Handler (qt)

/select

— common

q
obama; mexico

fq

sort

start, rows
0 10

fl

df

Raw Query Parameters
key1=val1&key2=val2

wt
json

☒ indent

☐ debugQuery

☐ dismax

☐ edismax

☐ hl

☐ facet

☐ spatial

☐ spellcheck

Execute Query

http://localhost:8080/discovery/select?q=obama%3B+mexico&wt=json&indent=true

```
{
  "responseHeader": {
    "status": 0,
    "QTime": 1,
    "params": {
      "indent": "true",
      "q": "obama; mexico",
      "_": "1425618506488",
      "wt": "json"
    }
  },
  "response": {
    "numFound": 0,
    "start": 0,
    "docs": []
  }
}
```

Fig. 2.1. Query: obama; mexico

The screenshot shows the Apache Solr Admin UI at `localhost:8080/#/discovery/query`. The left sidebar contains navigation links: Dashboard, Logging, Core Admin, Java Properties, Thread Dump, and a dropdown menu with Overview, Analysis, Dataimport, Documents, Files, Ping, Plugins / Stats, Query (selected), Replication, and Schema Browser.

The main panel is titled "Request-Handler (qt)" and shows the following fields:

- `/select`
- `common`
- `q`: chess
- `fq`
- `sort`
- `start, rows`: 0 to 10
- `fl`
- `df`
- `Raw Query Parameters`: key1=val1&key2=val2
- `wt`: json
- ☒ Indent
- ☐ debugQuery
- ☐ dismax
- ☐ edismax
- ☐ hl
- ☐ facet
- ☐ spatial
- ☐ spellcheck
- Execute Query

The right panel shows the JSON response for the query `http://localhost:8080/discovery/select?q=chess&wt=json&indent=true`:

```
{
  "responseHeader": {
    "status": 0,
    "QTime": 1,
    "params": {
      "indent": "true",
      "q": "chess",
      "_": "1425618782289",
      "wt": "json"
    }
  },
  "response": {
    "numFound": 1,
    "start": 0,
    "docs": [
      {
        "source_file_s": "MARCMerge20150306045430864244.warc@240340",
        "url": "http://www.chessgames.com/perl/chessgame?gid=1783846",
        "host": "www.chessgames.com",
        "domain": "chessgames.com",
        "public_suffix": "com",
        "server": [
          "Apache/2.2.8 (Ubuntu) PHP/5.2.4-2ubuntu5.27 with Suhosin-Patch mod_ssl/2.2.8 OpenSSL/0.9.8g mod_"
        ],
        "content_type_served": "text/html",
        "content_length": 21112,
        "id": "sha1:4KLRQWSQ00Y0V6V6GCPJQ0ZANIWSY2765/q54ftEBMphY+BVJSFabuFw==",
        "hash": [
          "sha1:4KLRQWSQ00Y0V6V6GCPJQ0ZANIWSY2765"
        ],
        "crawl_date": "2015-03-03T18:39:20Z",
        "crawl_year": "2015",
        "wayback_date": "20150303183920",
        "content": [
          "150fMembers · Prefs · Collections · Openings · Endgames · Sacrifices · History · Search Kibitzing · Kib_"
        ],
        "content_text_length": 3480,
        "content_type": [
          "text/html"
        ],
        "title": "Michael Adams vs David Baramidze (2015)",
        "headers": "chess chess games spiders chesski chessk BWM"
```

Fig. 2.2. Query: chess

The screenshot shows the Apache Solr Admin UI at `localhost:8080/#/discovery/query`. The left sidebar contains navigation links: Dashboard, Logging, Core Admin, Java Properties, Thread Dump, and a dropdown menu with Overview, Analysis, Dataimport, Documents, Files, Ping, Plugins / Stats, Query (selected), Replication, and Schema Browser.

The main panel is titled "Request-Handler (qt)" and shows the following fields:

- `/select`
- `q`: star wars; game of thrones
- `fq`
- `sort`
- `start, rows`: 0, 10
- `fl`
- `df`
- `Raw Query Parameters`: key1=val1&key2=val2
- `wt`: json
- ☒ indent
- ☐ debugQuery
- ☐ dismax
- ☐ edismax
- ☐ hl
- ☐ facet
- ☐ spatial
- ☐ spellcheck
- Execute Query

The right pane displays the JSON response for the query:

```
{
  "responseHeader": {
    "status": 0,
    "QTime": 1,
    "params": {
      "indent": "true",
      "q": "star wars; game of thrones",
      "wt": "json"
    }
  },
  "response": {
    "numFound": 1,
    "start": 0,
    "docs": [
      {
        "source_file_s": "WARCMerge20150306045430864244.warc@4328346",
        "url": "http://mehmetada.viralgalleries.me/star-wars-vs-game-of-thrones",
        "host": "mehmetada.viralgalleries.me",
        "domain": "viralgalleries.me",
        "public_suffix": "me",
        "content_type_served": "text/html; charset=UTF-8",
        "server": {
          "cloudflare-nginx"
        },
        "content_length": 114385,
        "id": "sha1:RQPJ8BWSGXJ7T6EV0QSKATY7IDW4KLD3/u560awdwqm41R9qdVt2zA==",
        "hash": [
          "sha1:RQPJ8BWSGXJ7T6EV0QSKATY7IDW4KLD3"
        ],
        "crawl_date": "2015-03-03T18:28:34Z",
        "crawl_year": "2015",
        "wayback_date": "20150303182834",
        "content": [
          "3bcavirallistfunnycelebrityastrologymoviesfoodtvNext ▶Star Wars vs Game of ThronesShareTweet#1Leia"
        ],
        "content_text_length": 3297,
        "content_type": [
          "text/html"
        ],
        "title": "Star Wars vs Game of Thrones",
        "content_encoding": "UTF-8",
        "content_ffb": "33626361",
        "content_first_bytes": "33 62 63 61 0d 0a 3c 21 44 4f 43 54 59 50 45 20 68 74 6d 6c 20 50 55 42 4c 49",
        "content_type_droid": "application/xhtml+xml; version=1.0",
        "links_hosts": [
          "imqur.com",

```

Fig. 2.3. Query: star wars; game of thrones

- From the previous query I toggled the case for the search parameter and the ranking of the results is modified as shown below.

The screenshot shows the Apache Solr Query interface in a web browser. The address bar indicates the URL is `localhost:8080/#/discovery/query`. The interface is divided into three main sections:

- Left Sidebar:** Contains navigation links for Dashboard, Logging, Core Admin, Java Properties, Thread Dump, and a dropdown menu for 'discovery'. Below this are links for Overview, Analysis, Dataimport, Documents, Files, Ping, Plugins / Stats, Query (highlighted), Replication, and Schema Browser.
- Central Query Form:**
 - Request-Handler (qt):** Set to `/select`.
 - q:** The search query is `StAr wArS`, with a red squiggly line indicating a spelling correction.
 - fq:** Empty.
 - sort:** Empty.
 - start, rows:** `0` to `10`.
 - fl:** Empty.
 - df:** Empty.
 - Raw Query Parameters:** `key1=val1&key2=val2`.
 - wt:** Set to `json`.
 - Options:** ☐ indent, ☐ debugQuery, ☐ dismax, ☐ edismax, ☐ hl, ☐ facet, ☐ spatial, ☒ spellcheck, ☐ spellcheck.build, ☐ spellcheck.reload.
- Right-Hand Pane:** Displays the JSON response for the query. The URL shown is `http://localhost:8080/discovery/select?q=StAr+wArS&wt=json&spellcheck=true`. The response is a JSON object with the following structure:


```
{
  "responseHeader": {
    "status": 0,
    "QTime": 1,
    "params": {
      "spellcheck": "true",
      "q": "StAr wArS",
      "_": "1425619690623",
      "wt": "json"
    }
  },
  "response": {
    "numFound": 3,
    "start": 0,
    "docs": [
      {
        "source_file_s": "WARCMerge20150306045430864244.warc@2956497",
        "url": "http://flavorwire.com/7p=502916",
        "url_type": "slashpage",
        "host": "flavorwire.com",
        "domain": "flavorwire.com",
        "public_suffix": "com",
        "server": [
          "nginx"
        ],
        "content_type_served": "text/html; charset=UTF-8",
        "content_length": 79207,
        "id": "sha1:200HP2EUKOFRTS3UGJ3QQVPSDTUKMHGM/1o2L02L9va0pn1PLgbpA4w==",
        "hash": [
          "sha1:200HP2EUKOFRTS3UGJ3QQVPSDTUKMHGM"
        ],
        "crawl_date": "2015-03-03T18:31:42Z",
        "crawl_year": "2015",
        "wayback_date": "20150303183142",
        "content": [
          "9c4FRSSTwitterFacebookFlavorwireSearchLocationMenuHot TopicsBroad CityFifty Shades of Grey"
        ],
        "content_text_length": 3609,
        "content_type": [
          "text/html"
        ]
      }
    ]
  }
}
```

Fig. 2.4. Query: StAr wArS

References

1. R bar plots. <http://www.theanalysisfactor.com/r-11-bar-charts/>.
2. SOLR front ends. <https://github.com/ukwa/webarchive-discovery/wiki/Front-ends>.
3. User guide for pywb. <https://github.com/ikreymer/pywb>.
4. WAIL. <http://matkelly.com/wail/>.
5. WARC and ARC indexing and discovery tools. <https://github.com/ukwa/webarchive-discovery>.
6. WARC merge. <https://github.com/maturban/WARCMerge>.