# Who and What Links to the Internet Archive

Yasmin AlNoamany, Ahmed AlSum, Michele C. Weigle, and Michael L. Nelson

Old Dominion University, Department of Computer Science
Norfolk VA 23529, USA
`{yasmin,aalsum,mweigle,mln}@cs.odu.edu`

**Abstract.** The Internet Archive's (IA) Wayback Machine is the largest and oldest public web archive and has become a significant repository of our recent history and cultural heritage. Despite its importance, there has been little research about how it is discovered and used. Based on web access logs, we analyze what users are looking for, why they come to IA, where they come from, and how pages link to IA. We find that users request English pages the most, followed by the European languages. Most human users come to web archives because they do not find the requested pages on the live web. About 65% of the requested archived pages no longer exist on the live web. We find that more than 82% of human sessions connect to the Wayback Machine via referrals from other web sites, while only 15% of robots have referrers. Most of the links (86%) from websites are to individual archived pages at specific points in time, and of those 83% no longer exist on the live web.

**Keywords:** Web Archiving, Web Server Logs, Web Usage Mining, Language Detection

## 1 Introduction

A variety of research has been conducted for studying web archives in order to answer questions related to user needs and to present web archive data to users [12,5]. However, no previous work has been carried out to answer these questions: What content languages are web archive users looking for? Why do users come to web archives? Where do web archive users come from? Who links to web archives? How do sites link to web archives? Do sites link deeply to specific archived pages or link to the repository? Why do sites link to the past?

The Internet Archive [11] is the first web archiving initiative attempting global scope and currently holds over 240 billion web pages with archives as far back as 1996 [8]. It allows traveling back in time for traversing archived versions of web pages through the Wayback Machine [18]. This paper provides a study of the requests of web archive users, both humans and robots, to gain insight into what users look for, in the context of the language of the requested pages, through an analysis of the server logs of the Internet Archives' Wayback Machine. We also provide an analysis of referring pages of human users to investigate how humans discover the Wayback Machine, why the referrers link to web archives, and how they link to web archives.

We found that users of Internet Archive's Wayback Machine request English pages the most, followed by several European languages. We also found that most human users come to the Wayback Machine via links or direct address presumably because they did not find the requested pages on the live web. Of the requested archived pages, 65% do not currently exist on the live web. From analyzing the referrers, we found that more than 82% of human sessions have referrers, while only 15% of robot sessions have referrers. We also found that 86% of the referrers are deep links to archived pages.

## 2    Related Work

To the best of our knowledge, no prior study has analyzed where web archive users come from nor what they look for in terms of the linguistic context. Furthermore, the usage of web archives in general has not been widely studied. The characterization of search behavior and the information needs of web archive users have been studied by Costa et al. [4,5] based on quantitative analysis of the Portuguese Web Archive (PWA) search logs. In a previous study [1], we provided the first analysis of user access to a large web archive. We discovered four basic access patterns for web archives through analysis of web server logs from the Internet Archive's Wayback Machine. In the study, we applied heuristics for robot detection after data filtering and found that robot sessions outnumber human sessions 10:1. Robots outnumber humans in terms of raw, unfiltered requests 5:4, and 4:1 in terms of megabytes transferred.

Many studies have investigated what is missing from digital libraries and web archives, in addition to the effect of this on the satisfaction of users' needs and expectations [17,3,22,16]. In [17], the Internet Archive's coverage of the web was investigated. The results showed an unintentional international bias through uneven representation of different countries in the archive. Carmel et al. [3] suggest a tool to dynamically analyze the query logs of the digital library system, identify the missing content queries, and then direct the system to obtain the missing data. We investigate what is missing through an analysis of requests with an HTTP 404 status in the Wayback Machine web server logs.
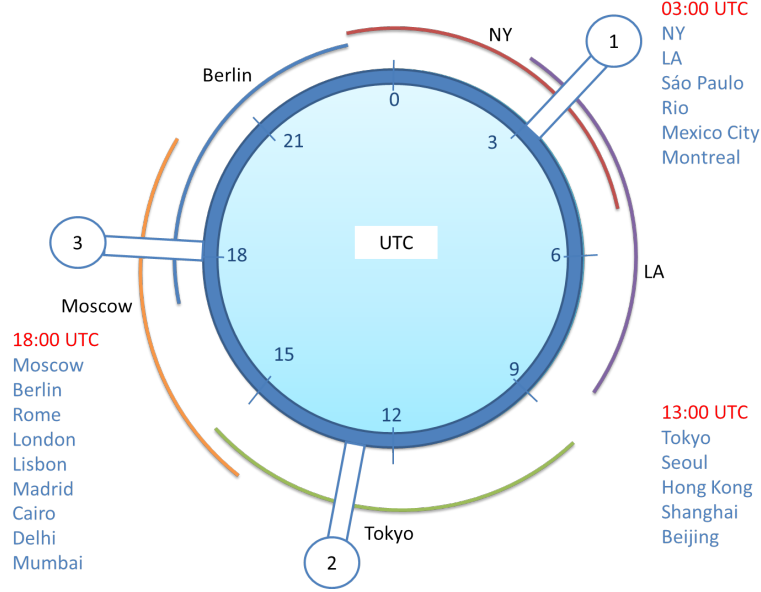
### Memento Terminology

In this section, we explain the terminology we adopt in the rest of the paper. Memento [20] is an HTTP protocol extension which enables time travel on the web by linking the current resources with their prior state. Memento defines the following terms:

- URI-R identifies the original resource. It is the resource as it used to appear on the live web. A URI-R may have 0 or more mementos (URI-Ms).
- URI-M identifies an archived snapshot of the URI-R at a specific datetime, which is called Memento-Datetime, e.g., URI-$M_i$= URI-R@$t_i$.

– URI-T identifies a TimeMap, a resource that provides a list of mementos (URI-Ms) for a URI-R with their Memento-Datetimes, e.g., $URI-T(URI-R) = \{URI-M_1, URI-M_2, ..., URI-M_n\}$.

Although we use Memento terminology, the logs we analyze are from the Internet Archive's Wayback Machine and not the Memento API.



**Fig. 1.** The dataset of 6M HTTP requests is constructed from slices of 2M each from 03:00, 13:00, and 18:00 UTC on February 2, 2012.

## 3  Methodology

We use the Internet Archive's Wayback Machine server logs in our analysis. We constructed our sample by combining three different slices of 2M records each (covering approximately 30 minutes) at times 03:00, 13:00, and 18:00 UTC on February 2, 2012, for a total dataset of 6M records. Because we are checking the language of the content accessed by the web archive users, we cover the peak time of Internet traffic periods for several countries with different language speakers to avoid biasing the results. According to many studies, the hours between 6 p.m. to 12 a.m. are considered to be peak times for Internet traffic [14,6,21]. We picked samples from the log file that were representative of the peak time for several cities around the world, as shown in Figure 1.

Table 1 contains the features of the sample. The features, from left to right, are the percentage of requests that used the GET method, were for embedded

| GET | Embedded | Null Ref | 2xx | 3xx | 4xx | 5xx | Humans | Robots |
|-----|----------|----------|-----|-----|-----|-----|--------|--------|
| 98.7% | 42.9% | 46.6% | 33.1% | 51.4% | 12.0% | 3.5% | 1.5% | 18.8% |

**Table 1.** Data set statistics based on 6M requests. Note that the last two columns are the percentage of humans and robots remaining after cleaning (removing the irrelevant requests to the analysis, e.g., embedded resources).

resources of web pages (such as images and CSS files), had null referrers (i.e., they do not identify a URI that links to a page at the Internet Archive), were successful requests (2xx status code), were redirections (3xx status code), were client errors (4xx status code), were server errors (5xx status code), remained from human requests after cleaning (removing the irrelevant requests to the analysis, e.g., embedded resources), and remained from robot requests after cleaning. The characteristics are consistent with our previous analysis of web archives [1].

Preparing the Wayback access logs for usage mining starts with transforming the raw log file into server sessions through web log preprocessing (data cleaning, user identification, and session identification) [13]. A session is the group of consecutive requests performed by a user [10]. We apply the same methodology as in our previous work for preprocessing the logs and web robot detection [1].

## 4   What do Wayback Machine Users Look for?

In this section, we give insight into what web archive users look for in terms of the content language of requested pages. We used the language detection library created by Shuyo [15] for detecting the language.

### 4.1   Archived Web Pages

**Distribution of Languages Used in the Wayback Machine** We extracted the successful requests (HTTP 200 status code) from humans and robots to detect the language distributions for the content of the requested pages. These successful requests represent 93.1% (85,909 out of 92,204) of all human requests and 56.7% (639,684 out of 1,127,204) of all robots requests. The request can be for a URI-T or a URI-M. For the URI-Ts, which represent 13% of human requested pages and 80.8% of robot requested pages, we estimated the language by using the most recent URI-M from the TimeMap. We identified 52 different languages from the successful requests. The left two columns of Table 2 show the top 10 languages which accounted for 94.8% of human and 93.4% of robot requests. For both human and robot users, English contributes the most to the successful requests, reflecting the high web archive penetration rate in English speaking countries. Japanese is the second most frequent language with 5.5% for humans, but Russian is the second most frequent language for robots at 7.0%. We also notice that despite of the existence of web archives in Europe, the requests to the IA from speakers of European languages contribute 13% of the top 10 list for human requested pages and 18.5% of the top 10 list for the robot requests.

| URI-Ms with HTTP 200 | | | | URI-Rs with HTTP 404 | | | |
|---|---|---|---|---|---|---|---|
| Language | Humans | Language | Robots | Language | Humans | Language | Robots |
| English | 71.7% | English | 72.4% | English | 66.9% | English | 62.2% |
| Japanese | 5.5% | Russian | 7.0% | Russian | 7.9% | Russian | 11.1% |
| German | 3.6% | German | 3.1% | German | 5.4% | German | 3.8% |
| Vietnamese | 2.9% | Spanish | 1.9% | Japanese | 5.1% | Indonesian | 3.1% |
| Russian | 2.3% | French | 1.8% | Spanish | 2.5% | Polish | 2.5% |
| Portuguese | 2.1% | Vietnamese | 1.7% | Polish | 2.3% | Vietnamese | 2.2% |
| French | 2.1% | Japanese | 1.5% | Romanian | 1.6% | Spanish | 2.0% |
| Spanish | 1.9% | Polish | 1.5% | French | 1.2% | Thai | 1.9% |
| Bengali | 1.8% | Portuguese | 1.3% | Italian | 0.8% | French | 1.8% |
| Italian | 0.9% | Thai | 1.1% | Portuguese | 0.7% | Dutch | 1.1% |

**Table 2.** The top 10 languages for URI-Ms with HTTP 200 (on the left) and for the URI-Rs of unarchived requested pages (on the right).

**Existence on the Live Web** From all 85,909 successful human requests, we checked the existence of the 40,791 unique URI-Rs on the live web. The robots generated 639,684 successful requests, in which there are 331,573 unique URI-Rs whose existence on the live web were also checked. We also checked the pages that give "soft 404s", which return HTTP 200, but do not actually exist, based on the algorithm in [2]. Table 3 contains the results of checking the status of the web pages on the live web.

| | Found in Archive | | Unarchived | |
|---|---|---|---|---|
| | Humans | Robots | Humans | Robots |
| **URI-Rs available on live web** | 36.4% | 62.5% | 25.4% | 33.2% |
| **URI-Rs missing from live web** | 63.6% | 37.5% | 74.6% | 66.8% |
| **Uniq. URI-Rs** | **40,791** | **331,573** | **2,441** | **209,384** |

**Table 3.** The existence of the requested archived pages on the live web. Available represents the requests which ultimately return "HTTP 200", while missing represents the requests that return HTTP 4xx, HTTP 5xx, HTTP 3xx to others except 200, timeouts, and soft 404s.

We believe humans access the Wayback Machine because they do not find web pages on the live web. Table 3 shows that for the requested pages that were found in the archive (returned HTTP 200 status), the percentage of the available pages on the live web for human requests is 36.4%. On the other hand, the percentage of the available pages on the live web for robot requests is 62.5%.

### 4.2   Unarchived Web Pages

Of the 6M requests in our sample, 12% returned HTTP 404 status, as shown in Table 1. Not all of these are actually unarchived; approximately 2% of the unique

URI-Rs are malformed (e.g., http://http://cnn.com) and were removed. We used the remaining valid URI-Rs (209,348 robots and 2,441 humans) to detect content language, check live web status, and check existence in other archives.

**Existence on the Live Web** The current state of the requested URI-Rs that had HTTP 404 status was determined by testing their existence on the live web. Of the URI-Rs that were not found in the Wayback Machine, 66.8% of those requested by robots and 74.6% of those requested by humans do not exist on the live web. To compensate for transient errors we repeated the requests several times for a week before declaring a URI-R non-existent.

**Distribution of the Content-Language for Unarchived Web Pages** We detected the content language of available URI-Rs on the live web, which represent 25.4% (620 out of 2,441) of the unique URI-Rs for humans and 33.2% (69,510 out of 209,384) for robots. The total number of requested URI-Rs is 227,450 for robots and 1,578 for humans. The two rightmost columns of Table 2 have the results for robots and humans separately. For the web pages that were not archived in IA's Wayback Machine, English is the most requested language with 66.9% of the human-requested web pages and 62.2% of the robot-requested web pages. The top 10 languages compromised 94.5% of all the content-language of the requested pages. European languages made up 22.5% of the human-requested pages and 22.4% of the robot-requested pages.

| Web Archive | Archive Web Site | #URI-R | #URI-M |
|---|---|---|---|
| Internet Archive (2013) | web.archive.org | 56,503 | 1,657,264 |
| The National Archives | webarchive.nationalarchives.gov.uk | 787 | 15,354 |
| ArchiefWeb | www.archiefweb.eu | 47 | 18,347 |
| Archive-It | archive-it.org | 41 | 4,682 |
| UK Web Archive | www.webarchive.org.uk | 38 | 12,277 |
| Library of Congress | webarchive.loc.gov | 35 | 1,092 |
| WebCite | webcitation.org | 29 | 1,104 |

**Table 4.** The number of the found URI-Rs and the corresponding URI-Ms of the missing pages (211,825 unique URI-Rs) on the web archives.

**Existence in Other Web Archives** We checked the 211,825 unarchived pages for existence in other archives at the time of the experiment. The existence in the web archives was tested by querying Memento proxies and aggregator [19]. For completeness and fairness, we also included the results from IA's Wayback Machine in March 2013. This resulted in 56,503 out of 211,825 URI-Rs that were unarchived in Feb. 2012 now being available in the archive. Table 4 contains the number of URI-Rs found in the web archives and the number of covered URI-Ms.

The Internet Archive has the most coverage at the time of experiment as they have increased their repository recently [8].

## 5    Where do Wayback Machine Users Come From?

We used the referrer field, which contains the web page that links to the resource, for the logs in our sample to determine how people discover the Wayback Machine. In terms of sessions, 84.8% of robot sessions do not have referrers while only 18.1% of human sessions do not have referrers (i.e., they reached the Wayback Machine by a link in an email, direct address, or direct bookmark). An empty referral field is a strong indicator of a robot.

| Web Site | Percentage | Description |
|---|---:|---|
| en.wikipedia.org | 12.9% | Wikipedia |
| archive.org | 11.9% | IA Home Page |
| reddit.com | 10.2% | Social News Web Site |
| google.TLD | 9.9% | Search Engine |
| info-poland.buffalo.edu | 1.5% | Polish Studies |
| de.wikipedia.org | 1.4% | Wikipedia |
| cracked.com | 1.2% | Humor Site |
| snopes.com | 1.1% | Urban Legends Reference Pages |
| facebook.com | 0.9% | Social Media |
| crochetpatterncentral.com | 0.9% | Crocheting Hobbies |

**Table 5.** The top 10 referrers.

In this section, we provide a detailed analysis of the referrer field of human users to gain insight into who links to the Wayback Machine and how they link to it. Robots are not included in the analysis of referrers because the majority of robots do not have referrers and if they do, we do not necessarily trust their values.

### 5.1    Who Links to Wayback Machine?

The percentage of human sessions with referrers is 81.9%. We eliminated the sessions that were referred by a URI-M or URI-T because they started prior to our sample. Of the sessions that started with an external referrer, 9.6% came from Google. The users who came from the home page of the IA contributed to 11.9% of the sessions with referrers. That means that many people start with the IA to access the Wayback Machine.

**Top Referrers** Table 5 contains the top 10 referrers that link to IA's Wayback Machine. The list of top 10 referrers represents 51.9% of all the referrers. As the table shows, en.wikipedia.org outnumbers all other sites including the

search engine and the home page of Internet Archive (archive.org). Note that
"google.TLD" represents Google search and 24 other pages from Google (e.g.,
http://www.google.com/about/company/history.html). Since the majority are
from Google search, we describe it as search engine. Facebook also appears as a
top referrer, which indicates that many people share links to the past.

| TLD | | .com | .org | .net | .jp | .ru | .de | .edu | .to | .uk | .info |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Percentage | | 45.4% | 33.9% | 8.4% | 1.8% | 1.4% | 1.4% | 1.1% | 0.7% | 0.6% | 0.5% |

**Table 6.** The top 10 TLDs of the referrers.

| ccTLD | | .com | .uk | .de | .ca | .jp | .pl | .nl | .ru | .fr | .br |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Percentage | | 56.7% | 6.0% | 5.3% | 4.8% | 3.7% | 2.2% | 1.9% | 1.7% | 1.5% | 1.4% |

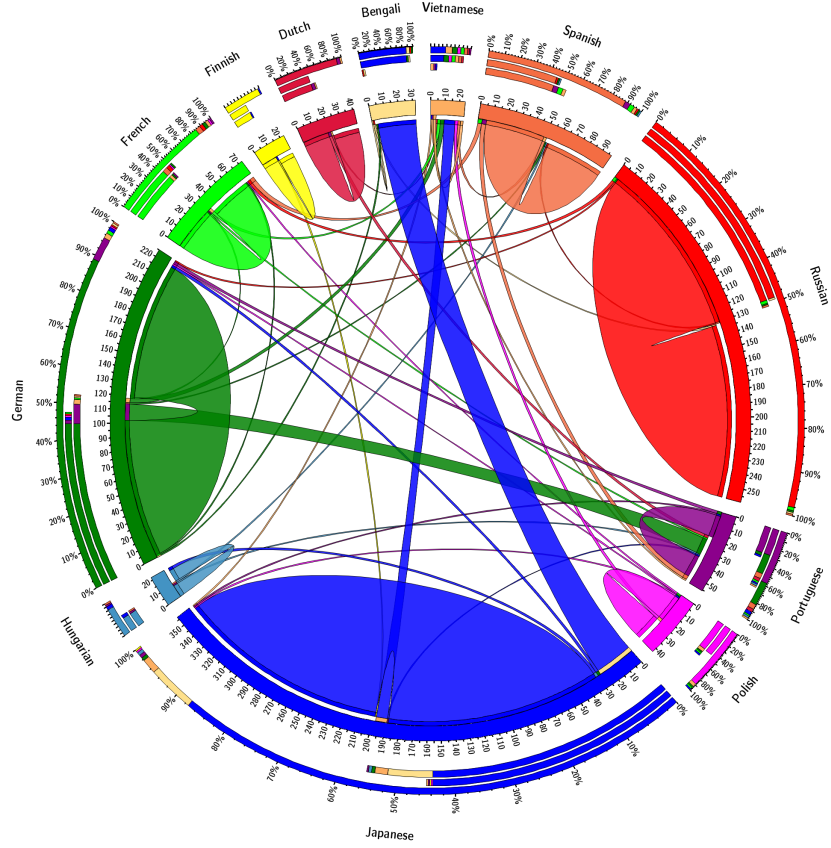**Table 7.** The top 10 ccTLDs of Google search referrers.

**Classification of Referrers** Table 6 presents the distribution of Top Level
Domains (TLD) for the URIs that link to the IA's Wayback Machine (only the
top 10 are shown). It can be noticed that most of the connections are from the
.com, .org, .net, .jp, .edu, and .ru domains. Despite of the existence of many web
archives in Europe, there are many European domains linking to the IA, such
as .ru (Russia), .de (Germany), .fr (France), and .it (Italy). Note that .to is the
TLD for a Russian language site (http://lurkmore.to/).

For the referrers from Google search, we extracted the country code top-level
domain (ccTLD) of the URIs to discover the countries of the users who came
to the Wayback Machine through the search engine. The results are shown in
Table 7. English-speaking countries are in the lead, followed by the European
language countries.

### 5.2   Inter-linking Between Languages

From the analysis of the content languages of the referrers and the archived
pages which have been linked by the referrers, English represents 80.7% of the
referrers' content languages and 80.2% of all referred pages. English referrers
link to English archived pages 92% of the time. A small percentage of English
referrers link to pages in other languages. The top 5 languages that English pages
link to are (in decreasing order) Portuguese, Vietnamese, French, and German.
Figure 2 contains a directed weighted graph, which is created using Circos [9], for
the relationship between the languages of the referrers and referees. We exclude
English from the graph to be able to analyze the rest of the languages and see
what they are linking to. For a particular language, the length of the outer arc
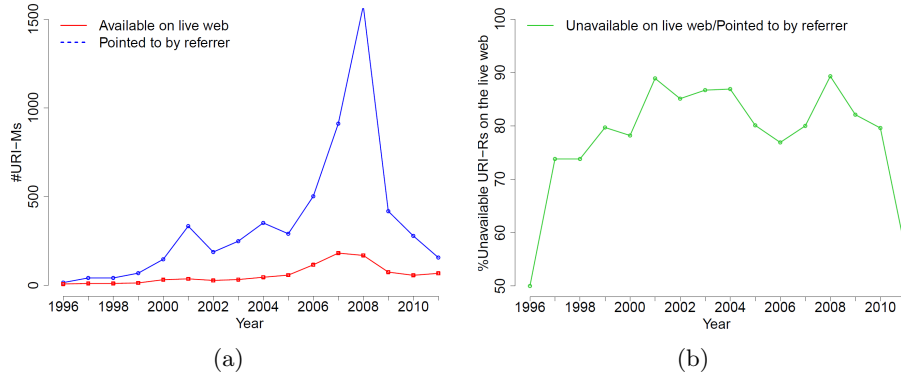
**Fig. 2.** Most languages self-link, with the notable exceptions of Japanese → {Bengali, Vietnamese} and German → Portuguese.

represents the sum of the number of referrer pages and the number of referee pages in that language. Moving toward the center, the next arc represents the percentage of referees, and the third arc represents the percentage of referrers in that language. For example, links to Japanese archived web pages denote 46% (160 out of 357) of all the Japanese language pages for referrers and archived pages all together. The inner circle shows the relationships between languages of the referrers and referees. Ribbons of different widths connect the languages. The direction is represented by a gap between the line and the incoming language (referrer language). For example, there are 30 links from Japanese (ja) pages to Bengali (bn) pages, which are shown as a fairly broad blue line. The languages where the relative number of referrer and referee together is less than 20 have been excluded to remove noise from the graph.

The figure shows that the languages are mainly linking to themselves with a few inter-language links. Though, recall that we have excluded English from the

**Fig. 3.** (a) The temporal distribution of URI-Ms pointed to by the referrers and the number of relative URI-Rs of these URI-Ms that are currently available on the live web. (b) The percentage of unavailable URI-Rs of these URI-Ms on the live web.

figure. Many of the top ranked languages of human-requested pages appear in the top ranked list of referrers, such as Japanese, German, Russian, Spanish, French, Polish, Dutch, Bengali, etc. It is surprising to find many European referrers to IA's Wayback Machine in spite of the existence of European web archives.

### 5.3   How do Web Pages Link to the Wayback Machine?

We found that 86.4% of the web pages that link to the Wayback Machine are pointing to mementos, which means they link to web pages at a specific time. There are 12.8% of web pages that point to TimeMaps. The percentage of web pages that point to the repository (e.g., http://web.archive.org) is 0.8%. Google search links to the top level URI, because Google does not crawl the archive based on the robots.txt exclusion protocol.

**Temporal Distribution of the Referred URI-Ms**  Figure 3(a) shows the total number of mementos which were pointed to by the referrers, grouped by the year of their Memento-Datetime. There is a significant bias toward 2008, then 2007, and then a bias against the more distant past. We found 14 URI-Ms all from a single web site that link to a datetime in 2099. We assume that the referrer wants to redirect the site's visitors to the most recent copy of the linked web page.

**Why do Web Sites Link to the Wayback Machine?**  The nature of the web is ephemeral, and the expected lifetime of a web pages is short [7]. So, web archives are important to webmasters and third parties for preserving and saving many web sites. Figure 3(b) clarifies that most people link to the Wayback Machine because they did not find the pages on the live web. The figure shows

that for most of the years, more than 70% of the referred pages on the archive no longer exist on the live web. About 83% of all referred-to URI-Rs do not currently exist on the live web.

## 6 Future Work and Conclusions

We plan to extend our analysis for investigating if the destination of users affects the session length and the behavior of web archive users. Furthermore, we will investigate the behavior of robots in web archives more and contrast it with the behavior of robots on the live web to distinguish their behaviors.

From the analysis of Internet Archive's Wayback Machine server logs, we conclude that most humans come to the Wayback Machine to find missing pages from the live web. The percentage of the requested archived pages which currently do not exist on the live web is 65%. We provided analysis for the distributions of languages to gain insight about what users look for. We found that English is the most used language on the Wayback Machine, followed by many European languages. European languages represent about 22% of the web pages that were not found on the Wayback Machine, for both human and robot requests. The large percentage of European languages among the unarchived pages can be a good indicator for archival demand for European web pages. We also provided analysis for the human referrers to discover where Wayback Machine users come from. We discovered that wikipedia is the most frequent referrer of pages to IA's Wayback Machine. From analyzing the TLDs of the referrers, we found many European domains (.ru, .de, .fr, etc.) in the top list of the referrers. English represents 80.2% of the referrer languages, followed by European languages. We found that the languages are linking mainly to themselves and to English. We also found that 86% of the referrer web pages link deeply to mementos. More than 82% of the links to these mementos are because their corresponding URI-Rs do not exist on the live web.

## 7 Acknowledgment

## References

1. AlNoamany, Y., Weigle, M.C., Nelson, M.L.: Access Patterns for Robots and Humans in Web Archives. In: Proceedings of the 13th ACM/IEEE-CS Joint Conference on Digital Libraries. JCDL '13 (July 2013)
2. Bar-Yossef, Z., Broder, A.Z., Kumar, R., Tomkins, A.: Sic transit gloria telae: towards an understanding of the web's decay. In: Proceedings of the 13th International Conference on World Wide Web. WWW '04, ACM (2004) 328–337

3. Carmel, D., Yom-Tov, E., Roitman, H.: Enhancing digital libraries using missing content analysis. In: Proceedings of the 8th ACM/IEEE-CS Joint Conference on Digital Libraries. JCDL '08, ACM (2008) 1–10

4. Costa, M., J. Silva, M.: Characterizing Search Behavior in Web Archives. In: Proceedings of Temporal Web Analytics Workshop. TWAW (2011)

5. Costa, M., Silva, M.J.: Understanding the Information Needs of Web Archive Users. In: Proc. of the 10th International Web Archiving Workshop. (Sept 2010)

6. Fukuda, K., Cho, K., Esaki, H.: The impact of residential broadband traffic on Japanese ISP backbones. SIGCOMM Comput. Commun. Rev. **35**(1) (January 2005)

7. Harrison, T.L., Nelson, M.L.: Just-In-Time Recovery of Missing Web Pages. In: Proceedings of the 17th Conference on Hypertext and Hypermedia. HYPERTEXT '06, ACM (2006) 145–156

8. Kahle, B.: Wayback Machine: Now with 240,000,000,000 URLs. `http://blog.archive.org/2013/01/09/updated-wayback/` (January 2013)

9. Krzywinski, M.I., Schein, J.E., Birol, I., Connors, J., Gascoyne, R., Horsman, D., Jones, S.J., Marra, M.A.: Circos: An information aesthetic for comparative genomics. Genome Research (2009)

10. Markov, Z., Larose, D.T.: Data Mining the Web: Uncovering Patterns in Web Content, Structure, and Usage. John Wiley & Sons, Inc. (2007)

11. Negulescu, K.C.: Web Archiving @ the Internet Archive. Presentation at the 2010 Digital Preservation Partners Meeting, `http://1.usa.gov/XSjDG8` (2010)

12. Padia, K., AlNoamany, Y., Weigle, M.C.: Visualizing Digital Collections at Archive-It. In: Proceedings of the 12th ACM/IEEE-CS Joint Conference on Digital Libraries. JCDL '12, ACM (2012) 15–18

13. Reddy, K.S., Varma, G.P.S., Babu, I.R.: Preprocessing the Web Server Logs  An illustrative approach for effective usage mining. ACM SIGSOFT Software Engineering Notes **37**(3) (May 2012) 1–5

14. Reisinger, D.: Netflix gobbles a third of peak Internet traffic in North America. CNET, `http://goo.gl/2cVPg` (2012)

15. Shuyo, N.: Language Detection Library for Java. `http://code.google.com/p/language-detection/` (2012)

16. Silva, A.J.C., Gonçalves, M.A., Laender, A.H.F., Modesto, M.A.B., Cristo, M., Ziviani, N.: Finding what is missing from a digital library: A case study in the computer science field. Inf. Process. Manage. **45**(3) (May 2009) 380–391

17. Thelwall, M., Vaughan, L.: A fair history of the web? examining country balance in the internet archive. Library & Information Science Research **26**(2) (2004)

18. Tofel, B.: Wayback for Accessing Web Archives. In: Proceedings of International Web Archiving Workshop. IWAW (2007)

19. Van de Sompel, H., Nelson, M.L., Sanderson, R.: HTTP framework for time-based access to resource states – Memento. `https://datatracker.ietf.org/doc/draft-vandesompel-memento/` (2012)

20. Van de Sompel, H., Nelson, M.L., Sanderson, R., Balakireva, L.L., Ainsworth, S., Shankar, H.: Memento: Time Travel for the Web. Technical Report arXiv:0911.1112 (2009)

21. Wasserman, T.: Netflix takes up 32.7% of Internet bandwidth. Marshable, `http://goo.gl/2FtWa` (2011)

22. Zhuang, Z., Wagle, R., Giles, C.: What's there and what's not?: focused crawling for missing documents in digital libraries. In: Proceedings of the 5th ACM/IEEE-CS Joint Conference on Digital Libraries. JCDL '05 (2005) 301–310