

Analyzing Web archives through Topic and Event Focused Sub-Collections

Gerhard Gossen
L3S Research Institute
Leibniz Universität Hannover
Appelstraße 9a
30167 Hannover, Germany
gossen@L3S.de

Elena Demidova
Web and Internet Science
Group, Electronics &
Computer Science
University of Southampton
e.demidova@soton.ac.uk

Thomas Risse
L3S Research Institute
Leibniz Universität Hannover
Appelstraße 9a
30167 Hannover, Germany
risse@L3S.de

ABSTRACT

Web archives capture the history of the Web and are therefore an important source to study how societal developments have been reflected on the Web. However, the large size of Web archives and their temporal nature pose many challenges to researchers interested in working with these collections. In this work, we describe the challenges of working with Web archives and propose the research methodology of extracting and studying sub-collections of the archive focused on specific topics and events. We discuss the opportunities and challenges of this approach and suggest a framework for creating sub-collections.

CCS Concepts

• **Applied computing** → **Digital libraries and archives**;
• **Information systems** → **Web mining**; *Document filtering*;

Keywords

Web archive; sub-collection; topics; events

1. INTRODUCTION

Web archives such as the Internet Archive¹ or the archives collected by national libraries allow researchers in Web Science and the Digital Humanities to look back at the past of the Web and trace its development over time. These archives are created by regularly crawling the Web (in the case of the Internet Archive) or selected subsets (typically national sub-domains) to create snapshots of Web sites at different points in time. Researchers can look up any of the crawled versions to look back at specific points in the past or compare different versions.

An important challenge when using Web archives is the access to the collected data. As an example, the Internet

¹<http://www.archive.org>

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

WebSci '16, May 22 - 25, 2016, Hannover, Germany

© 2016 Copyright held by the owner/author(s). Publication rights licensed to ACM. ISBN 978-1-4503-4208-7/16/05...\$15.00

DOI: <http://dx.doi.org/10.1145/2908131.2908175>

Archive has a size of several petabytes over a time span of 20 years. A researcher working with such an archive needs efficient and effective tools to scope relevant documents for further research. In contrast to typical use cases where only individual pages are considered (e.g. in legal disputes or to provide persistent links) or the entire archive is analyzed using automatic methods (e.g. using text mining) [12], many research questions in Web Science and the Digital Humanities require an analysis of documents related to specific topics and events. The analysis is often performed manually, therefore the documents to be analyzed in detail have to be carefully selected. We call a collection of documents related to a specific topic or event a **topic and event focused sub-collection** of the archive.

Current tools do not support the researcher enough in creating such a sub-collection (cf. for example [7]). Current approaches use *browsing* and *searching* as access methods for Web archives. *Browsing* is done by entering URLs in a Web interface such as the Wayback Machine [8] and navigating the archived pages using hyperlinks. This requires that the URLs of relevant pages or at least pages linking to them are known in advance. As the entire process is done manually, the researcher will typically have to try many hyperlinks that are not available in the archive. The researcher also has to be aware of temporal drift while navigating the archive [1]. Temporal drift occurs because the linking and linked page were usually crawled at different times and therefore each navigation step moves the analyzed point in time. Iterated navigation can even lead to page versions that were crawled outside of the relevant time window. The browsing approach is therefore inefficient and error-prone.

An alternative approach is to use keyword *search* over the entire archive. Here the user only needs to enter keywords into a search interface and can see all pages matching the query. Many search interfaces also provide faceted browsing, where the search can be narrowed down further by e.g. the crawl time, the document content type or the domain name. It is much easier for users to get started using this approach, as they do not have to know about relevant documents in advance. They can also combine this approach with the browsing approach by starting the navigation from a search result document. Search has however many technical challenges. First of all, the archive needs to create and keep up to date a full text index of its entire content and execute queries over this index. Because of the typically large size of the archives this necessitates distributed indexing and retrieval architectures [5] that require many computational

Scope	Type	Description
URL	list of URLs	documents that need to be in the sub-collection
domain	list of domains	domains that the sub-collection should be restricted to
time	time interval	relevant timeframe
keywords	list of keywords	descriptive keywords for the sub-collection topic
event/entity	list of knowledgebase references	entries in a knowledgebase such as FreeBase [2] that are the topic of the sub-collection
size	number of documents	target size of the sub-collection

Table 1: Exemplary scopes used in a sub-collection specification. This list is not exhaustive.

resources as well as a lot of technical expertise to set up and maintain. An additional issue is the ranking of documents matching a query. On the one hand, the archive contains many snapshots of the same document, which can negatively impact traditional ranking methods that use for example link-based measures [4]. On the other hand, the information need of Web archive users is different than in standard search applications because it is usually focused around specific time periods and therefore requires different ranking measures [13]. Therefore the searching of Web archive using current tools still requires a lot of manual effort from the user to tune queries and go through result lists.

We therefore propose an alternative approach for the access to Web archives through the automatic extraction of topic and event focused sub-collections. Such sub-collections contain documents from the archive that have been automatically classified as referring to a given topic or event. We may additionally pose the constraints that documents in the collection are connected through hyperlinks or that the temporal distance between any two documents in the collection is minimal. By extracting such sub-collections we provide researchers with relevant sets of documents, which can be further analyzed in their appropriate context.

In this work we define the concept of topic and event focused sub-collections, describe a framework for extracting them and discuss challenges.

2. WEB ARCHIVE SUB-COLLECTIONS

In this section we define the concept of topic and event focused sub-collections and describe several important variants of such sub-collections.

A **Web archive** is a collection of **Web document snapshots**. Web document snapshot refers to the content retrieved from a given URL (**document URL**) at a given time (**crawl time**). In addition to the content, the archive typically also stores metadata about the snapshot such as the software used for retrieval, the HTTP headers or the document content type.

An **topic and event focused sub-collection** is a set of Web document snapshots, where each snapshot is available in the given Web archive. It is defined in a **sub-collection specification** that describes the scope of the sub-collection. The format of the sub-collection specification depends on the **extraction algorithm** used to create the sub-collection. Each type of algorithm defines a number of **scopes** that describe relevant documents. A list of exemplary scopes is given in Tab. 1. Scopes can be combined to narrow down the sub-collection. For example, a simple algorithm that supports the **URL** and **time** scopes can be used to extract all snapshots of a given URL in a specific time frame. A scope

does not need to be exact. For example, a topic scope can be implemented using a machine learning algorithm, that classifies a snapshot as relevant based e.g. on the similarity to a given set of topic keywords. In this case evaluation metrics like precision and recall can be used to analyze the quality of a given algorithm.

Given the large size of typical Web archives it is often neither feasible nor desired to find all snapshots matching a scope. Therefore an algorithm should have a high precision in matching the scopes but may have a lower recall. A good algorithm should however aim to find a representative sub-collection, i.e. one that has a similar diversity as the original archives in terms of domains, crawl times or types of sources.

We further distinguish between a **connected** and a **disconnected** sub-collection. A connected sub-collection needs to contain for any snapshot s contained in the sub-collection also at least one snapshot t for each document that is linked from s , if one is available in the archive. In contrast, a disconnected sub-collection can consist only of isolated snapshots. A connected sub-collection is needed to perform e.g. link graph analyses, whereas e.g. content-based analyses can also be performed on a disconnected sub-collection.

An additional distinction is between **snapshot** and **timeline** sub-collections. In a snapshot sub-collection, each document URL should occur only once, a timeline sub-collection should however have all snapshots of an in scope URL that are also in scope. A snapshot sub-collection is useful in synchronic analyses, where the researcher is looking at a specific point in time and does not want to deal with multiple versions of the same URL. In contrast, a timeline sub-collection is needed to perform diachronic analyses where we want to track a development over time.

3. SUB-COLLECTION EXTRACTION FRAMEWORK

In this section we will describe the framework, in which the sub-collection extraction process takes place and describe measures to evaluate extraction algorithms.

To create a Web archive sub-collection C the researcher first has to choose a base Web archive W and create a sub-collection specification CS that describe their collection need. Then they need to select an algorithm A that supports the scopes specified in CS and run it over the archive W , using the sub-collection specification CS as a parameter. The result of this process is the extracted sub-collection C . We expect that the extraction process will typically be iterative, i.e. that the research will create a modified specification CS' after analyzing the sub-collection C and create a new sub-collection C' , maybe even using a different extraction algorithm A' .

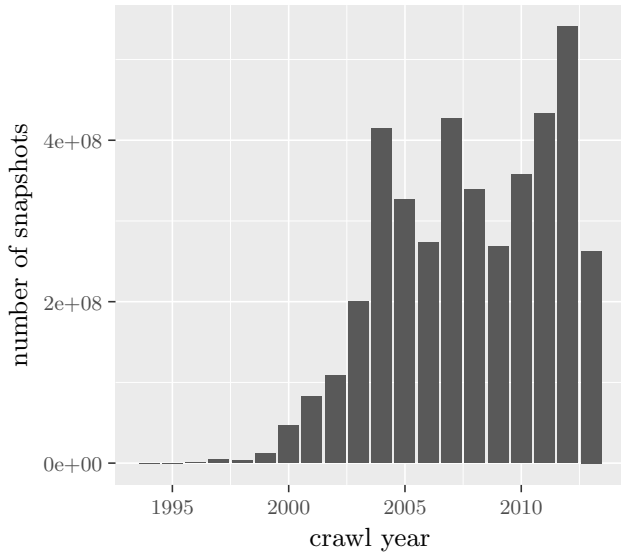


Figure 1: Number of snapshots by crawl time

To compare different extraction algorithms, we can use the following evaluation measures:

Precision As there is no metadata about what topics and events a snapshot in a Web archive is relevant for, the relevance calculation needs to be done using automated methods such as machine learning, e.g. Support Vector Machines (SVMs) for topic classification [9]. These methods will however mistake some irrelevant documents as relevant and vice versa. The precision measures the rate of such errors (cf. [11]):

$$\text{precision} := \frac{|\text{retrieved relevant snapshots}|}{|\text{retrieved snapshots}|}$$

Recall An extraction algorithm can simply iterate over the entire archive and evaluate the specified scopes on each snapshot. Given the large size of Web archives, this is often prohibitively expensive even if parallel processing facilities are available. Therefore it is desirable that the extraction algorithm can use indexes or heuristics to speed up the execution. The recall measures the fraction of relevant snapshots that were extracted from the archive (cf. [11]):

$$\text{recall} := \frac{|\text{retrieved relevant snapshots}|}{|\text{relevant snapshots}|}$$

In this way it quantifies the expected loss of using a more efficient algorithm.

Note that we also need to consider the precision of the relevance estimation method when computing the recall as it may perform differently on the selected subset, for example because the extraction algorithm will preferentially select pages from certain domains or having a certain link structure. In this case also the recall on the entire collection needs to be examined.

Diversity As described above, the goal of extracting sub-collections is to find a set of documents that help in

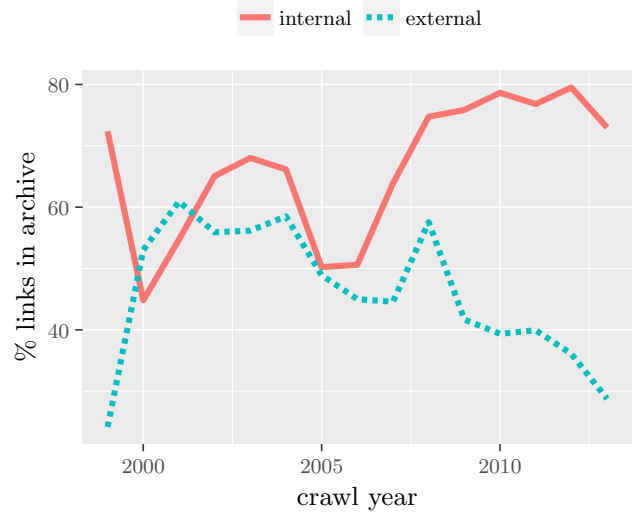


Figure 2: Average rate of outlinks that are contained in the archive. Internal links have the same domain name as the linking page, external links go to different domains.

answering a research question. While the collection needs to have a manageable size so that it can be analyzed, it also needs to be representative of the entire archive. This can be measured using diversity measures that describe how well different aspects of the given topic or event are represented [3].

Link completeness When we analyze the context of a snapshot, e.g. using link graph analysis, it is important to also have all relevant linked pages in the sub-collection. We measure the link completeness of a collection C as follows:

$$lc(C) = \sum_{s \in C} \frac{|\text{retrieved relevant outlinks of } s|}{|\text{relevant outlinks of } s|}$$

Temporal coherence Snapshots in a Web archive are typically crawled at different points in time, even if they refer to the same event. Additionally, a given URL may have been crawled several times in a relevant time frame, providing an algorithm with multiple snapshots to choose from. The selection of snapshots can however introduce errors in the downstream analyses when selecting snapshots that are from distant points in time. A way to reduce this risk is to optimize the selection of snapshot such that the time between any pair of snapshots is minimized. A similar measure is the *blur* of a Web archive which also considers the expected number of changes to retrieved pages [6].

Run time To allow for a fast iteration of refined sub-collection specifications, it is important that extraction algorithms can produce their results fast. As operations on Web archives are often executed in parallel on large clusters, typical run time measures such as the elapsed time in

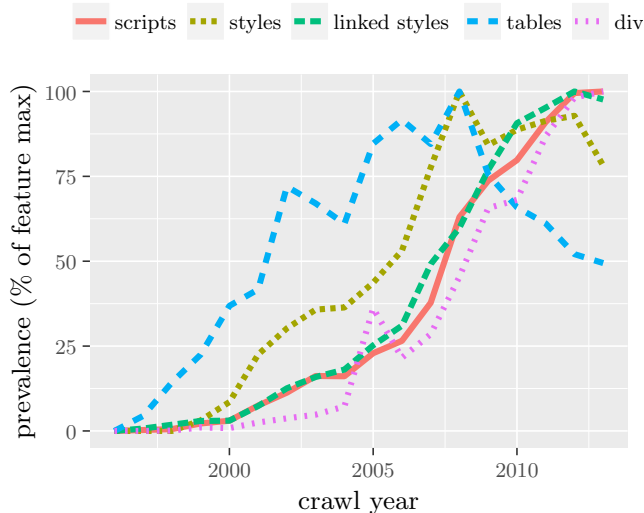


Figure 3: Average number occurrences of HTML tags per document. Each value is normalized to its maximum value.

seconds are less useful because they are heavily influenced by the size and current load of the cluster. A better evaluation metric is instead the number of disk accesses to retrieve and evaluate snapshots, as this is typically the dominant cost in the extraction process.

4. CHALLENGES

In this section we discuss several challenges when working with Web archives, especially when extracting sub-collections. We illustrate the challenges using data extracted from the German Web archive, a collection of all Web pages from the .de top level domain crawled by the Internet Archive between 1994 and 2013. All values except the total number of snapshots (Fig. 1) were calculated using a random sample of 40K snapshots.

Temporal Scope The number of snapshots per year can have strong fluctuations (see Fig. 1). In general there are more snapshots for recent years, which is consistent with the general trend of a growing web. The exact number of snapshots per year can however vary due to different crawl strategies, intermittent errors and other factors.

In the context of Web archive sub-collections this means that a temporal scope can only be used effectively for some time periods, whereas e.g. in the first years of the archive’s time span it may exclude too many documents. This also means that diversification may be necessary to avoid that snapshots from sparse time periods get lost while more active time periods are over-represented.

Archive Completeness Most archives do not have complete snapshots of the Web due to limited resources,

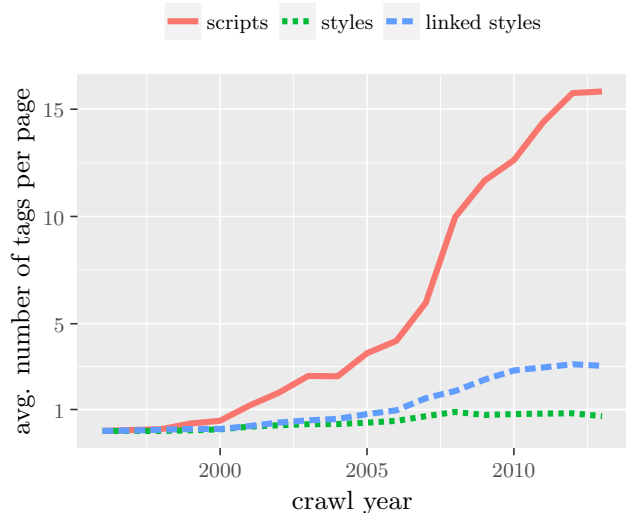


Figure 4: Average number occurrences of HTML tags per document.

legal restrictions and technical challenges [10] that restrict the collection of the archives. Fig. 2 gives estimates of rate of missing content. For each snapshot in our sample we extracted the outgoing links and tried to retrieve them from the archive. We see that about 50-80% of all links within the same domain (internal links) can be retrieved, whereas for external links this rate is much lower at 20-60%. A possible reason for the lower ratio for external links is the restriction of the archive to the .de domain. This is however a realistic scenario, as many current Web archives are run by national libraries and have similar restrictions to Web sites from specific countries or top level domains.

For the sub-collection extraction this means that there is an inherent upper bound on the achievable link completeness. It also suggests that optimizing for link completeness may bias the sub-collection towards sites with many relevant internal links and away from hub pages with many relevant external links, as the former are likely to be more complete.

Content Diversity The content in our archive spans 20 years, which means that it reflects many developments of Web technology. Fig. 3 shows the relative frequency of representative HTML tags over time. Each curve has been normalized such that it reflects the prevalence of the tag in comparison to its peak value. We can easily see the decline of table-based layouts in the 2000s and the growth of layout techniques using div tags styled by CSS style sheets. Similarly, Fig. 4 shows the absolute average number of scripts and style sheets per page. We see after 2005 a dramatic increase in the number scripts and to a lesser extent of linked style sheets. The former may be a reflection of the increased use of advertising networks and tracking services.

Similarly, we can look at the number of links per page over time (Fig. 5). Whereas the number of external

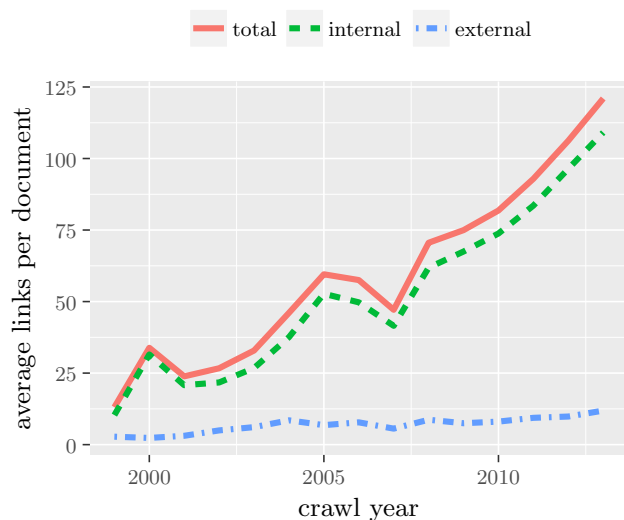


Figure 5: Average number of outlinks per document. Internal links have the same domain name as the linking page, external links go to different domains.

links stays relatively stable, the number of internal links increases continuously. This may be due to the maturing of Web sites that accumulate content over time, but it may also reflect the changed Web environment where e.g. search engine optimization (SEO) through specific link strategies becomes more common.

For the extraction of sub-collections this means that our algorithms must be adaptive to different types of Web content: Over time the format of Web pages has changed dramatically, therefore also the position of relevance cues on the pages may have changed. Furthermore, the value of links has changed over time, such that we have to be more selective when selecting links on more recent pages.

5. CONCLUSION

In this work we have presented a new approach for the access to Web archives through the automatic extraction of topic and event focused sub-collections. In contrast to existing approaches, sub-collection decrease the amount of manual effort required to find a reasonably-sized collection of documents for further research, while increasing the value of the collection through better extraction of link graphs and the avoidance of temporal drift. We have defined the problem of extracting sub-collections and have described several typical extensions to the problem. Additionally, we have described the framework in which this approach is applied and have shown several evaluation metrics to compare algorithms for this problem. Based on data from a real-world Web archive we have demonstrated several issues for algo-

rithms trying to solve this problem and have discussed approaches for dealing with them. We hope that this work sparks interest in the extraction and use of topic and event focused sub-collections. In future work, we will present algorithms that create sub-collections following this framework.

Acknowledgments

This work was partially funded by the European Research Council under ALEXANDRIA (ERC 339233) and the European Commission under SoBigData (RIA 654024) and H2020-MSCA-ITN-2014 WDAqua (grant agreement 64279).

6. REFERENCES

- [1] S. Ainsworth and M. L. Nelson. Evaluating sliding and sticky target policies by measuring temporal drift in acyclic walks through a web archive. In *Joint Conference on Digital Libraries*, JCDL '13, pages 39–48, 2013.
- [2] K. Bollacker, C. Evans, P. Paritosh, T. Sturge, and J. Taylor. Freebase: a collaboratively created graph database for structuring human knowledge. In *SIGMOD International Conference on Management of data*, SIGMOD '08, pages 1247–1250. ACM, 2008.
- [3] C. L. Clarke, M. Kolla, G. V. Cormack, O. Vechtomova, A. Ashkan, S. Büttcher, and I. MacKinnon. Novelty and diversity in information retrieval evaluation. In *SIGIR Conference on Research and Development in Information Retrieval*, SIGIR '08, pages 659–666, 2008.
- [4] M. Costa, F. Couto, and M. Silva. Learning temporal-dependent ranking models. In *SIGIR Conference on Research and Development in Information Retrieval*, SIGIR '14, 2014.
- [5] M. Costa, D. Gomes, F. M. Couto, and M. J. Silva. A survey of web archive search architectures. In *Temporal Web Analytics Workshop 2013*, Rio de Janeiro, Brazil, 2013.
- [6] D. Denev, A. Mazeika, M. Spaniol, and G. Weikum. The sharc framework for data quality in web archiving. *The VLDB Journal*, 20(2):183–207, Apr. 2011.
- [7] R. Deswarte. Revealing british euroscepticism in the uk web domain and archive case study, July 2015. Available at <http://sas-space.sas.ac.uk/id/eprint/6103>.
- [8] International Internet Preservation Consortium (IIPC). OpenWayback. <http://netpreserve.org/openwayback>, 2016. [Online; accessed 2016-01-15].
- [9] T. Joachims. Text categorization with Support Vector Machines: Learning with many relevant features. In *European Conference on Machine Learning*, ECML '98, pages 137–142, 1998.
- [10] M. Kelly, M. L. Nelson, and M. C. Weigle. The archival acid test: Evaluating archive performance on advanced HTML and JavaScript. In *Joint Conference on Digital Libraries*, JCDL '14, pages 25–28, Piscataway, NJ, USA, 2014. IEEE Press.
- [11] C. D. Manning, P. Raghavan, and H. Schütze. *Introduction to Information Retrieval*. Cambridge University Press, New York, NY, USA, 2008.
- [12] E. Reynolds. Web archiving use cases. Technical report, International Internet Preservation Committee, 2013. Available at <http://netpreserve.org/resources/web-archiving-use-cases-0>.
- [13] G. Weikum, N. Ntarmos, M. Spaniol, P. Triantafyllou, A. A. Benczúr, S. Kirkpatrick, P. Rigaux, and M. Williamson. Longitudinal analytics on web archive data: It's about time! In *Conference on Innovative Data Systems Research (CIDR '11)*, pages 199–202, 2011.