

Article

ARCOMEM Crawling Architecture

Vassilis Plachouras^{1,*}, Florent Carpentier², Muhammad Faheem³, Julien Masanès²,
Thomas Risse⁴, Pierre Senellart³, Patrick Siehndel⁴, Yannis Stavrakas¹

¹ Institute for the Management of Information Systems, ATHENA Research and Innovation Center, Artemidos 6 & Epidavrou, Maroussi 15125, Greece

² Internet Memory Foundation, 45 ter rue de la Révolution, 93100 Montreuil, France

³ CNRS LTCI, Institut Mines-Télécom, Télécom ParisTech, 46 rue Barrault, 75634 Paris Cedex 13, France

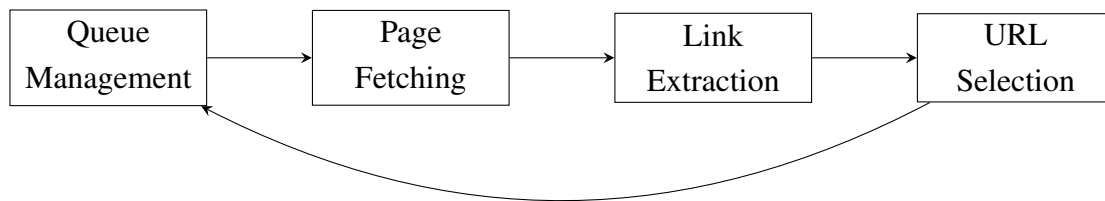
⁴ L3S Research Center, University of Hannover, Appelstr. 9a, 30167 Hannover, Germany

* Author to whom correspondence should be addressed; vplachouras@acm.org, Tel. +302106875413, Fax +302106856804.

Version July 31, 2014 submitted to *FutureInternet*. Typeset by \LaTeX using class file *mdpi.cls*

1 **Abstract:** The World Wide Web is the largest information repository available today.
2 However, this information is very volatile and Web archiving is essential to preserve it for the
3 future. Existing approaches to Web archiving are based on simple definitions of the scope
4 of Web pages to crawl and are limited to basic interactions with Web servers. The aim of
5 the ARCOMEM project is to overcome these limitations and to provide flexible, adaptive
6 and intelligent content acquisition, relying on social media to create topical Web archives.
7 In this article, we focus on ARCOMEM's crawling architecture. We introduce the overall
8 architecture and we describe its modules, such as the online analysis module which computes
9 a priority for the Web pages to be crawled, and the Application-Aware Helper which takes
10 into account the type of Web sites and applications to extract structure from crawled content.
11 We also describe a large-scale distributed crawler that has been developed, as well as the
12 modifications we have implemented to adapt Heritrix, an open source crawler, to the needs
13 of the project. Our experimental results from real crawls show that ARCOMEM's crawling
14 architecture is effective in acquiring focused information about a topic and leveraging the
15 information from social media.

16 **Keywords:** Web archiving; crawling architecture; content acquisition

Figure 1. Traditional processing chain of a Web crawler

17 1. Introduction

18 The World Wide Web is the largest information repository. But this information is very *volatile*: the
19 typical half-life of content referenced by URLs is of a few years [1]; this trend is even aggravated in
20 social media, where social networking APIs sometimes only extend to a week's worth of content [2].
21 Web archiving [3] deals with the collection, enrichment, curation, and preservation of today's volatile
22 Web content in an archive that remains accessible to tomorrow's historians. Different strategies for Web
23 archiving exist: bulk harvesting, selective harvesting and combinations of both. Bulk harvesting aims at
24 capturing snapshots of entire domains. In contrast selective harvesting is much more focused, e.g., on
25 an event or a person. Combined strategies include less frequent domain snapshots complemented with
26 regular selective crawls. In the following we will focus on the technical aspects of selective crawls.

27 Selective crawls require a lot of manual work for the crawl preparation, crawler control, and quality
28 assurance. On the technical level, current-day archiving crawlers, such as Internet Archive's Heritrix [4],
29 crawl the Web in a conceptually simple manner (See Figure 1). They start from a *seed* list of URLs
30 (typically provided by a Web archivist) to be stored in a queue. Web pages are then fetched from this
31 queue one after the other, stored as is in the archive, and further links are extracted from them. If newly
32 extracted links point to URLs that are in the scope of archiving tasks (usually given by a list or regular
33 expressions of URLs to consider), they are added to the queue. This process ends after a specified time
34 or when there is no interesting URL left to crawl. Due to this simple way of crawling, bulk domain
35 crawls are well supported while selective crawls necessitate additional manual work for the preparation
36 and quality assurance. It is the aim of the ARCOMEM¹ project to support the selective crawling on the
37 technical level by leveraging social media and semantics to build meaningful Web archives [5]. This
38 requires, in particular, a change of paradigm in how content is collected technically via Web crawling,
39 which is the topic of the present article.

40 This traditional processing chain of a Web crawler like Heritrix [4] has several major limitations:

- 41 • Only regular Web pages, accessible through hyperlinks and downloadable with an HTTP GET
42 request, are ever candidates for inclusion in the archive; this excludes other forms of valuable
43 Web information, such as that accessible through Web forms, social networking RESTful APIs, or
44 AJAX applications.
- 45 • Web pages are stored as is in the archive, and the granularity of the archive is that of a Web page.
46 Modern Web applications, however, often present individual blocks of information on different

¹ <http://www.arcomem.eu/>

47 parts of a Web page: think of the messages on a Web forum, or the different news items on a news
48 site. These individual *Web objects* can be of independent interest to archive users.

- 49 • The crawling process does not vary from one site to another. The crawler is blind to the kind of Web
50 application hosted by this Web site, or to the software (typically, a *content management system*)
51 that powers this Web application. This behavior might lead to resource loss in crawling irrelevant
52 information (e.g., login page, edition page in a wiki system) and prevents any optimization of the
53 crawling strategy within a Web site based on how the Web site is structured.
- 54 • The scope of a selective crawl is defined by a crude whitelist and blacklist of URL patterns; there
55 is no way to specify that relevant pages are those that are related to a given semantic entity (say, a
56 person) or that are heavily referenced from influential users in social networks.
- 57 • The notion of scope is binary: either a Web page is in the scope or it is not – on the other hand, it
58 is very natural for a Web archivist to consider various degrees of relevance for different pieces of
59 Web content; and ideally content should be crawled by decreasing degree of relevance.

60 The crawling architecture of ARCOMEM aims at solving these different issues by providing flexible,
61 adaptive, intelligent content acquisition. This is achieved by interfacing traditional Web crawlers such
62 as Heritrix with additional modules (complex resource fetching, Web-application-aware extraction and
63 crawling, online and offline analysis of content, prioritization), as well as by adapting the internals of
64 the crawlers when needed (typically for managing priorities of content relevance). The objective of this
65 article is to present an overview of this crawling architecture, and of its performance (both in terms of
66 efficiency and of quality of the archive obtained) on real-Web crawls. This article extends [6].

67 The remainder of this work is organized as follows. We first discuss in Section 2 the related work.
68 Then we present in Section 3 a high-level view of the ARCOMEM architecture, before reviewing
69 individual modules in Section 4. We present evaluation results that highlight the effectiveness of
70 ARCOMEM's crawling architecture in Section 5. Finally, we present our concluding remarks in
71 Section 6.

72 2. Related Work

73 While crawling appears to be a simple process, there are several associated challenges, especially
74 when the aim is to crawl a large number of Web pages [7], in order to create the index of a Web search
75 engine, or to archive them for future reference.

76 **Web Crawling.** Descriptions of early versions of Google's and Internet Archive's large-scale crawler
77 systems appeared in [8] and [9], respectively. However, one of the first detailed descriptions of a
78 scalable Web crawler is that of Mercator by Heydon and Najork [10], who provide information on the
79 various modules of the crawler and the design options. Najork and Heydon also describe a distributed
80 crawler based on Mercator in [11]. Shkapenyuk and Suel [12] introduce a distributed and robust crawler,
81 managing the failure of individual servers. Heritrix [13] is an archival-quality and modular open source
82 crawler, developed at the Internet Archive. In Section 4.4 we will describe how we have adapted Heritrix

83 in order to fit in ARCOMEM's crawling architecture. Boldi et al. [14] describe UBIcrawler, a distributed
84 Web crawler, implemented in Java, which operates in a decentralized way and uses consistent hashing
85 to partition the domains to crawl across the crawling servers. Lee et al. [15] describe the architecture
86 and main data structures of IRLBot, a crawler which implements DRUM (Disk Repository with Update
87 Management) for checking whether a URL has been seen previously. The use of DRUM allows IRLBot
88 to maintain a high crawling rate, even after crawling billions of Web pages.

89 As the Web evolves, and Web pages are created, modified, or deleted [16][17], effective crawling
90 approaches are needed to handle these changes. Cho and Garcia Molina [18] describe an incremental
91 crawler for optimizing the average freshness of crawled Web data. Olston and Pandey [19] describe
92 re-crawling strategies to optimize freshness based on the longevity of information on Web pages. Pandey
93 and Olston [20] also introduce a parameterized algorithm for monitoring Web resources for updates and
94 optimizing timeliness or completeness depending on application-specific requirements.

95 **Focused and Deep-Web Crawling.** Focused or topical crawlers [21] provide an effective way to
96 balance the cost, coverage, and quality aspects of data collection from the Web [22], by selectively
97 crawling pages that are relevant to a set of topics, defined as a set of keywords [23], by example
98 documents mapped to a taxonomy of topics [24], or by ontologies [25][26]. Recent approaches also
99 address the crawling of information for specific geographical locations [27][28].

100 The main challenges in focused crawling relate to the prioritization of URLs not yet visited, which
101 may be based on similarity measures [23][25], hyperlink distance-based limits [29][30], or combinations
102 of text and hyperlink analysis with Latent Semantic Indexing (LSI) [31]. Machine learning approaches,
103 including naïve Bayes classifiers [24][32], Hidden Markov Models [33], reinforcement learning [34],
104 genetic algorithms [35], and neural networks [36], have also been applied to prioritize the unvisited
105 URLs.

106 Focused crawlers and crawlers in general can harvest data from the publicly indexable Web by
107 following hyperlinks between Web pages. However, there is a very large part of the Web that is hidden
108 behind HTML forms [37]. Such forms are easy to complete by human users. Automatic deep-Web
109 crawlers, however, need to complete HTML forms and retrieve results from the underlying databases.
110 Barbosa and Freire [38] develop mechanisms for generating simple keyword queries that cover the
111 underlying database through unstructured simple search forms. Madhavan et al. [39] handle structured
112 forms by automatically completing subsets of fields, aiming to obtain small coverage over many hidden
113 Web databases.

114 **Web Archiving.** Web archiving refers to the collection and long-term preservation of data available
115 on the Web [3]. Since archiving the whole Web is a very challenging task due to its size and dynamics,
116 there have been several national initiatives for preserving the Web of a country, based on full crawls in
117 Sweden [40] and on a selective collection of Web pages in the United Kingdom [41] and Australia [42].
118 The former approach aims at providing complete snapshots of a domain taken at regular intervals. A
119 drawback of this approach is the lack of knowledge about changes of Web pages between crawls and the
120 consistency of the collected data [43]. The latter approach results in higher quality collections restricted
121 only to selected Web sites. Denev et al [44] introduce a framework for assessing the quality of archives

122 and tune the crawling strategies to optimize quality with given resources. Gomes et al. [45] provide a
123 survey of Web archiving initiatives.

124 Focused crawlers, as described above, can be used for creating focused Web archives, by relying on
125 a selective content acquisition approach. The crawling process in the ARCOMEM architecture changes
126 the paradigm in how content is collected technically via Web crawling, by performing selective crawls
127 and also leveraging information found in online social media.

128 3. Crawling Architecture

129 The goal for the development of the ARCOMEM crawler architecture was to implement a socially
130 aware and semantic-driven preservation model [5]. This requires thorough analysis of the crawled Web
131 page and its components. These components of a Web page are called *Web objects* and can be the title,
132 a paragraph, an image, or a video. Since a thorough analysis of all Web objects is time-consuming, the
133 traditional way of Web crawling and archiving is no longer functioning. Therefore the ARCOMEM
134 crawl principle is to start with a *semantically enhanced crawl specification* that extends traditional
135 URL-based seed lists with semantic information about entities, topics or events. This crawl specification
136 is complemented by a small reference crawl to learn more about the crawl topic and intention of the
137 archivist. The combination of the original crawl specification with the extracted information from
138 the reference crawl is called the *Intelligent Crawl Definition (ICD)*. This specification, together with
139 relatively simple semantic and social signals, is used to guide a broad crawl that is followed by a thorough
140 analysis of the crawled content. Based on this analysis a semi-automatic selection of the content for the
141 final archive is carried out.

142 The translation of these steps into the ARCOMEM crawling architecture foresees the following
143 processing levels: the *crawler level*, the *online processing level*, the *offline processing level*, and the
144 *cross-crawl analysis* that revolve around the ARCOMEM database as depicted in Figure 2. Since the
145 focus of this article is the crawling and the online analysis we will focus on these levels in the rest of the
146 article and give only a brief overview on the other levels. More details about the other processing levels
147 and the whole architecture can be found in [5].

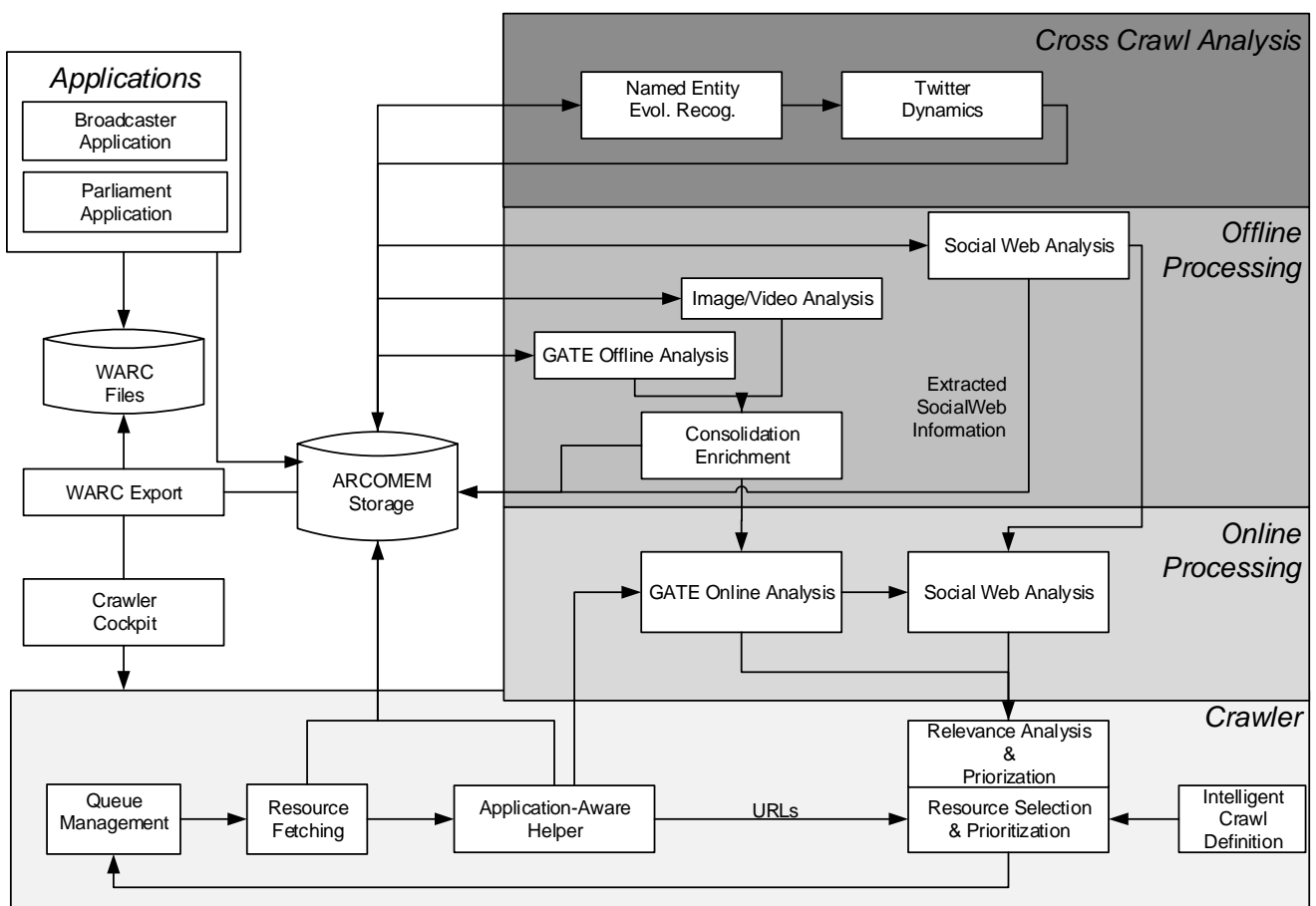
148 3.1. Crawling Level

149 At this level, the system decides and fetches the relevant Web objects as these are initially defined
150 by the archivists, and are later refined by both the archivists and the online processing modules. The
151 crawling level includes, besides the traditional crawler and its decision modules, some important data
152 cleaning, annotation, and extraction steps (we explain this in more detail in Section 4). The Web
153 objects (i.e., the important data objects existing in a page, excluding ads, code, etc.) are stored in the
154 ARCOMEM database together with the raw downloaded content.

155 3.2. Online Processing Level

156 The online processing is tightly connected with the crawling level. At this level a number of
157 semantic and social signals such as information about persons, locations, or social structure taken

Figure 2. Overall Architecture



158 from the intelligent crawl specification are used to prioritize the crawler processing queue. Due to
159 the near-real-time requirements, only time-efficient analysis can be performed, while complex analysis
160 tasks are moved to the offline phase. The logical separation between the online processing level and
161 the crawler level will allow the extension of existing crawlers at least with some functionalities of the
162 ARCOMEM technology. More details about the online analysis can be found in Section 4.2.

163 *3.3. Offline Processing Level*

164 At this level, most of the basic processing over the data takes place. The offline, fully-featured,
165 versions of the entity, topics, opinions, and events analysis (ETOE analysis) and the analysis of the social
166 contents operate over the cleansed data from the crawl that are stored in the ARCOMEM database. These
167 processing tools perform linguistic, machine learning and NLP methods in order to provide a rich set
168 of metadata annotations that are interlinked with the original data. In addition to processing of textual
169 content, multimedia content can also be analyzed and enriched with meta-information. The respective
170 annotations are stored back in the ARCOMEM database and are available for further processing and
171 information mining. After all the relevant processing has taken place, the Web pages to be archived and
172 preserved are selected in a semi-automatic way and transferred to the Web archive (in the form of WARC
173 files).

174 *3.4. Cross-Crawl Analysis Level*

175 Finally, a more advanced processing step takes places. It operates on collections of Web objects that
176 have been collected over time and can cover several crawls. Analysis implemented exemplary on this
177 level within the ARCOMEM system is used to recognize Named Entity Evolutions [46] and to analyze
178 the evolutions of associations between interesting terms and tweets (Twitter Dynamics) [47].

179 *3.5. Applications*

180 For the interaction with the crawler and exploration of the content a number of applications are used
181 around the ARCOMEM core system. The crawler cockpit is used to create the crawl specification, to
182 monitor the crawl activities, and to initiate the final export of crawled content to WARC files.

183 The end-user applications allow users to search archives by domain, time, and keywords.
184 Furthermore, browsing the archives via different facets such as topics, events, and entities, and
185 visualizing the sentiments of Social Web postings complement the end-user application. However, the
186 applications are not limited to the described examples. The ARCOMEM system is open to any kind of
187 application that wants to use it.

188 4. Module description

189 The modular crawling architecture introduced in Section 3 enables the integration of a wide range of
190 functionalities and technologies in the same system. In this section, we describe in detail the Application-
191 Aware Helper, the online analysis module, as well as the two crawlers that can be used to acquire content.

192 4.1. Application-Aware Helper

193 The goal of the application-aware helper (AAH) is to make the crawler aware of the particular kind
194 of Web application it is crawling, in order to adapt the crawling strategy accordingly. The presence of
195 the AAH in the crawling processing chain ensures Web content is crawled in an intelligent and adaptive
196 manner.

197 The AAH applies different crawling strategies for different types of social Web sites (Web forums,
198 blogs, social networks, photo networks, music networks, video networks, etc.), or for specific content
199 management system (vBulletin, WordPress, etc.). The AAH detects the Web application and Web page
200 type within the Web application before deciding which crawling strategy is best for the given Web
201 application.

202 More precisely, this module performs the following steps:

- 203 1. it detects the Web application type (general type, content management system, etc.);
- 204 2. it detects the Web page type (e.g., in a Web forum, if we are at the level of a list of forums, a list
205 of threads, or a list of messages);
- 206 3. it executes the relevant crawling actions; extracting Web objects (e.g., comments, posts) on the
207 one hand, and adding only relevant URLs to the crawlers' queue on the other hand.

208 The AAH is assisted by a *knowledge base*, which specifies how to detect a specific Web application
209 and which crawling actions should be applied. The knowledge base is written in a custom XML format,
210 so as to be easily shared, updated, and hopefully managed by non-programmers. The knowledge base
211 ensures that an appropriate crawling strategy is applied for a detected Web application. For instance,
212 the vBulletin Web forum CMS can be identified by searching for a reference to a specific script with
213 the detection pattern: `script[contains(@src,'vbulletin_core.js')`]. The AAH distinguishes
214 two main kinds of Web application *levels*: *intermediate* pages, such as lists of forums, lists of threads,
215 can only be associated with navigation actions; *terminal* pages, such as the individual posts in a forum
216 thread, can be associated with both navigation and extraction actions. For intelligent crawling, our
217 AAH needs not only to distinguish among Web application types, but among the different kinds of
218 Web pages that can be produced by a given Web application type. For example, the expression `//h2`
219 `[@class="forumtitle"]/a/@href` detects *intermediate* pages in vBulletin, whereas the expression
220 `//table[@class="post"]` identifies *terminal* pages. Once the application type and level is detected,
221 the system executes the relevant crawling actions. The crawling actions are of two types: *extraction*
222 *actions* point to the individual Web objects to be extracted from a given Web page (e.g., comments, blog
223 post); *navigation actions* point to the URLs to be added in a crawling queue. For instance, the *extraction*
224 *action* `//div[contains(@id,'post_message')`] extracts the post message.

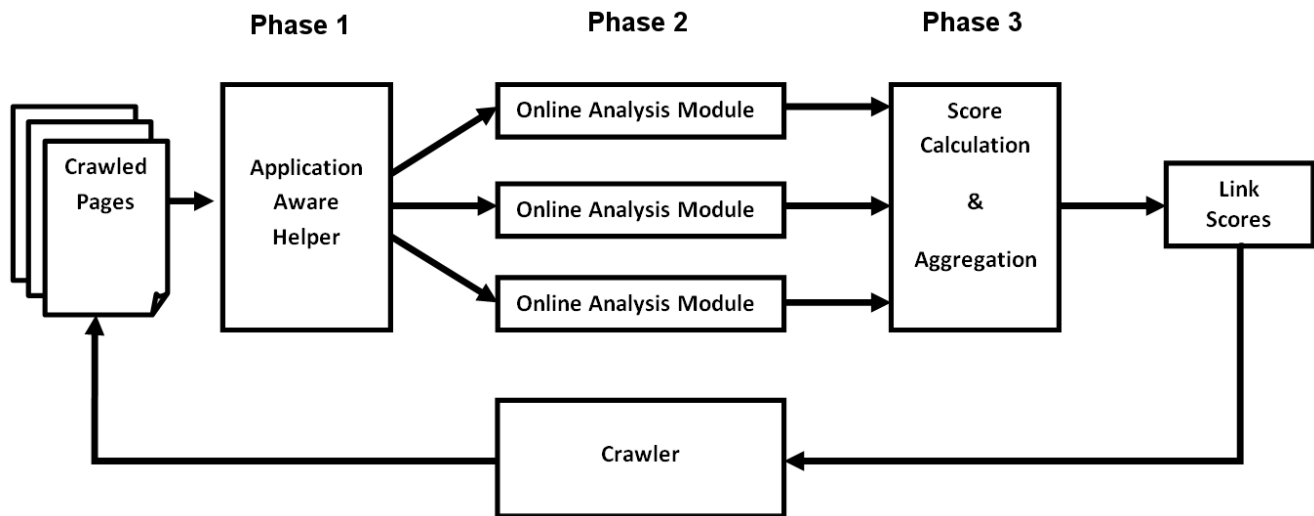
225 To exploit the AAH in Web-scale crawling, a *Web application adaptation* module has been integrated.
226 The AAH deals with both Web content changes and Web structure changes. When Web content changes,
227 the AAH simply updates recently crawled versions with the new one. However, the Web structure
228 changes (i.e., changes in Web site template) are more complicated to identify and adapt to. The AAH
229 deals with this challenge. Structural changes with respect to the knowledge base come from varying
230 versions of the CMS, or alternative templates proposed by CMSs or developed for a specific Web
231 application. Here, we assume that Web application detection patterns never fail. In our experiments,
232 we did not see any instance where a Web application was not successfully detected.

233 The AAH deals with two different cases of structure adaptation: first, when (part of) a Web application
234 has been crawled before, but recrawl of Web application fails after template changes; second, when a new
235 Web application has been detected successfully, but (some of) the existing actions are inapplicable. The
236 adaptation module for recrawls of Web application relearns appropriate crawling actions for each failed
237 crawlable object. In ARCOMEM, crawled Web pages with their Web objects and metadata are stored
238 in the form of RDF triples in the ARCOMEM database. Therefore, we have proposed an algorithm
239 which utilizes the ARCOMEM database and first detects structural changes for already crawled Web
240 applications by looking for the crawled content in the Web pages with crawling actions used during a
241 previous crawl. If the system fails to extract the content, then it means that the structure of the Web
242 application has changed. In the presence of structural changes, the algorithm detects the inappropriate
243 crawling actions and performs updates by aligning them according to structural changes.

244 In the case of a new Web application whose template is slightly different from the one present in
245 the knowledge base, the adaptation module cannot be applied on previously crawled content. Here, the
246 adaption is applied for two different scenarios: *Web application level detected* and *Web application not*
247 *detected*. We consider two classes of Web application levels: intermediate and terminal. The navigation
248 actions are applicable for the intermediate level (e.g., list of blog posts), whereas the terminal level
249 may require both navigation and extraction actions (e.g., individual post). When a Web application
250 level is detected, but (some) crawling actions fail, a set of relaxed expressions are generated by relaxing
251 predicates, and tag names. The candidate tag names are selected from the knowledge base as well as
252 from the existing page DOM tree. For example, if an expression `div[contains(@class,'post')]/`
253 `h2[@class='posttitle']` fails to extract the post title, and `div[contains(@class,'post')]` is the
254 detection pattern that worked, then we will try several relaxations of the second half of the expression,
255 for instance, replacing `@class` with `@id`, `'posttitle'` with `'posthead'`, `h2` with `div`, etc. We
256 favor relaxations that use parts from crawling actions in the knowledge base for other Web application
257 types of the same general category (e.g., bulletin boards). Any successful relaxed expression will be still
258 tested with a few more pages of the same Web application level.

259 When a Web application level is not detected then an appropriate crawling strategy cannot be initiated,
260 therefore the system first adapts the detection patterns. The idea here is the same as above: the system
261 first collects all the candidate attributes, tag names and values; and then creates all possible combinations
262 of relaxed expressions. For example, assume that the candidate set of attributes and values are: `@class`
263 `'post'`, `@id=forumlist`, `@class='bloglist'` with candidate set of tag names `article`, `div`, etc.
264 The set of relaxed expression will be generated by trying out each possible combination:

265 `// div[contains(@class,'post')]`

Figure 3. Interaction of online analysis modules

266 // div[contains(@id,'forumlist')]

267 // div[contains(@class,'bloglist')]

268 and similarly for other tag names. If the system detects the Web application with any relaxed expression
 269 then the system will apply the crawling actions adaptation as described above.

270 The AAH has reduced bandwidth, time, and storage (by requiring fewer HTTP requests for
 271 known Web applications, avoiding duplicates) using limited computational resources in the process.
 272 Application-aware crawling also helps adding semantic information to the ARCOMEM database. More
 273 details about the functioning and independent evaluation of the AAH are provided in [48,49].

274 4.2. Online Analysis

275 Within the online analysis, several modules analyze crawled Web objects in order to guide the crawler.
 276 The purpose of this process is to provide scores for detected URLs. These scores are used for guiding
 277 the crawler in order to obtain a focused crawl with respect to the intelligent crawl definition (ICD). The
 278 main modules used within the online analysis are the AAH, the GATE platform for text analysis [50],
 279 and a prioritization module. The actual online processing consists of three phases which are displayed in
 280 Figure 3; within Figure 2 these phases are related to the connections between the online processing and
 281 the crawler.

- 282 1. The AAH performs preprocessing steps on the crawled Web pages;
- 283 2. Online analysis modules run on relevant document parts;
- 284 3. The output of online analysis modules is aggregated and a score for each URL is provided.

285 A more detailed description of these steps is provided in the remainder of this section.

286 In the first phase we run the AAH on the Web page to detect regions of interest in the document
 287 and discard irrelevant parts. A detailed description of the functions provided by the AAH is given in

288 Section 4.1. The input document is split into one or more document parts. Each document part is
289 processed separately from now on.

290 In the second phase the online analysis modules are run on the content of the document part. Currently
291 we use textual analysis modules using GATE, a URL scoring module using URL patterns and a simple
292 spam link filter using a black list. Additional modules can be added easily. The textual analysis module
293 performs basic NLP pre-processing on the text and allows the extraction of relevant entities.

294 The version of GATE used within the Online Analysis is a lightweight version of GATE since the
295 performance and the processing time needs to be as fast as possible. The tasks carried out by the
296 GATE component comprise basic linguistic processing steps, language identification and Named Entity
297 Recognition (NER). In contrast to the use of the basic GATE functions which are needed to create
298 word vectors describing the crawled objects, the use of the NER module is optional. With the NER
299 module disabled the processing time of a Web document was reduced by about 70%. Based on these
300 observations we disabled the NER module in the online phase for most crawls and moved this part to
301 the offline analysis, where performance aspects are not as critical as during the online phase. Using the
302 extracted keywords and the given crawl specification we calculate a score based on the cosine similarity
303 of the term vectors. The matching is run at several granularities: whole document, paragraph around
304 anchor, and only anchor text. This allows us to boost link anchors that are closer to keyword or entity
305 matches.

306 Each analysis module can produce a score for the current document and one for each out-link. Some
307 analysis modules (e.g., the URL analyzers) omit the document score, while others can only provide
308 document scores (e.g., the text analysis). In the latter case the document score is propagated to each
309 out-link contained in the analyzed document.

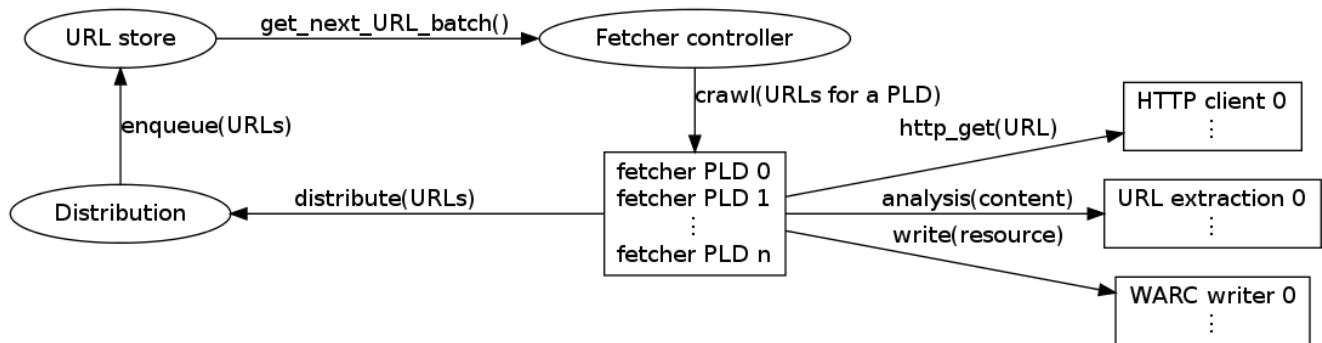
310 The final phase of the online analysis is the priority aggregation: The scores provided by the individual
311 analysis modules are aggregated into one final score for each out-link. Here we use a weighted average
312 over the individual scores using weights provided by the users.

313 4.3. Large-Scale Crawler

314 The large-scale crawler is a distributed crawler, implemented by the Internet Memory Foundation
315 (IMF). It retrieves content from the Web and stores it in WARC files, that can optionally be loaded into
316 an HBase repository. Its main initial aim was scalability: crawling at a fast rate from the start and slowing
317 down as little as possible as the amount of visited URLs grows to billions, all while observing politeness
318 conventions (rate regulation, robots.txt compliance, etc.). This objective is achieved by incorporating
319 recent developments in data structures and design options for crawlers [14][15]. It ran a crawl on 19
320 virtual machines with 8 cores and 32 GiB of RAM each for three weeks. The rate was kept over 2000
321 HTTP requests per second for the whole duration of the crawl for a total of close to four billion URLs
322 crawled.

323 The crawler does not require distributing a static node list to all cluster instances nor does it require
324 external utilities to copy lists of URLs as they get discovered. It also detects nodes joining or leaving the
325 cluster and changes the URL distribution mapping to account changes without any manual intervention.

Figure 4. Crawler architecture



326 Figure 4 depicts the main processes on each cluster node. The rectangles depict many processes with
 327 the same function. The ovals represent individual processes or subsystems made of many processes.

328 The fetcher controller is in charge of spawning fetchers, limited to spawning as many as is allowed
 329 by the configuration (mainly to respect memory constraints) and by the availability of URLs in the URL
 330 store. It asks the URL store for a batch of URLs all belonging to the same pay level domain (PLD,
 331 approximated by searching for the longest applicable public suffix and adding one more level), resolves
 332 the domain name to an IP address and ensures no other fetcher in the entire cluster is crawling this IP
 333 address. It then spawns a fetcher and passes it the URL batch. The fetcher gets the robots.txt file and
 334 starts crawling all the allowed URLs, respecting the required delay between each fetch.

335 For each resource, three main steps are performed:

- 336
- 337 • fetching (HTTP request);
 - 338 • analysing the document according to its type in search of new URLs. It may also run other analyses
 339 which may be useful at run time; for example, language identification;
 - 340 • writing the content plus extracted or derived information into a WARC file (depending on the
 341 configuration and filtering settings);
 - 342 • filtering according to the scope configuration before sending to the distribution module.

342 When a fetcher has processed all of its URLs, it exits and the fetcher controller will try to replace it
 343 with a new fetcher and a fresh batch of URLs.

344 The distribution module maintains a consistent hashing ring that reflects the current cluster topology.
 345 It forwards URLs to the appropriate node for them to be queued in the local URL store.

346 The URL distribution being based on the pay level domain, it is easy to guarantee that no more
 347 than one fetcher in the whole cluster will be crawling a specific host at any time. However, there is no
 348 guarantee that different pay level domains are not mapped to the same IP address. To ensure rate control
 349 for IP addresses, we use a global IP address registry.

350 The WARC files get copied asynchronously to a specific directory in a Hadoop file system (HDFS).
 351 A periodic import task will insert the content from the HDFS into HBase. This makes the crawler quite
 352 independent from the storage system. In particular, the crawler can continue to work without HBase for
 353 as long as it has available disk space.

354 To follow the numerous events occurring inside the crawler as tens of thousands of concurrent
355 processes run, a flexible system is necessary. We have implemented node-local filtering of events by
356 subsystem and severity, and centralized storage in a full text index that allows complex queries and
357 advanced graphical representations.

358 The IMF crawler can perform multiple crawls concurrently, supporting one URL store and a
359 configuration (scope functions, archival functions, etc.) for each concurrent crawl, while having a single
360 fetcher pool. This feature guarantees that politeness is respected across all crawls while allowing to
361 crawl concurrently as many domains as possible.

362 The latest developments are geared towards greater flexibility to add ease-of-use and ‘archive quality’
363 to the crawler’s scalability.

364 The large-scale crawler supports HTTPS, can stream large files, retries on server failures, detects
365 the real MIME type and language of documents, extracts many metadata from HTML pages (such as
366 outlinks with type, anchor text, etc.). It has a fast C implementation of a comprehensive and configurable
367 URL canonicalization. It provides advanced scoping functions that can be combined at will, allowing,
368 for instance, to make decisions based on the language of a page or the whole path that led to it. It also
369 employs a fully-fledged and extensible per-domain configuration framework with parameters including
370 budget, minimum and maximum delay between two fetches. Crawler fetchers subscribe to updates of
371 parameter values and use the new configuration immediately. It detects traps by analyzing URLs and
372 checking for similar content.

373 *4.4. Adaptive Heritrix*

374 In addition to the large-scale crawler developed by IMF, we have also investigated to what extent
375 Heritrix, a widely-used open source crawler [13], can be adapted to ARCOMEM’s crawling architecture.
376 Heritrix implements a typical centralized crawling process, where a URL is prioritized only when it is
377 added to the frontier. In order to adapt Heritrix to the needs of ARCOMEM, we have implemented
378 a frontier that supports updating the priorities of already scheduled URLs and receiving scored URLs
379 from external processes, possibly running on different servers. As a result, Heritrix can be used as
380 a fetching service for selective Web harvesting. Overall, we have extended Heritrix with a range of
381 functionalities, regarding the storing of crawled content, the extraction of anchor text with links, etc. All
382 the modifications are available as open source software in the releases of the ARCOMEM project. In
383 this article we focus on the two main features mentioned above.

384 The default frontier of Heritrix employs a Berkeley DB backed hash table for storing URLs, typically
385 grouped according to the domain or host they belong to. The key of a URL’s record in the frontier is
386 computed based on its domain, a flag indicating whether the URL should be crawled immediately, its
387 priority (or precedence in Heritrix terminology), and a counter, which increases for every URL that is
388 inserted in the frontier (Figure 5). The frontier implementation provides a *next* method for obtaining
389 the next URL to crawl from a given domain or host, but there is no method to update the priority of an
390 already scheduled URL.

391 To overcome this limitation, we have implemented a frontier that extends the default frontier of
392 Heritrix, adding a hash table that maps a URL already scheduled to the key with which it was scheduled.

Figure 5. The structure of the key corresponding to an entry of a URL scheduled for crawling in the default frontier of Heritrix

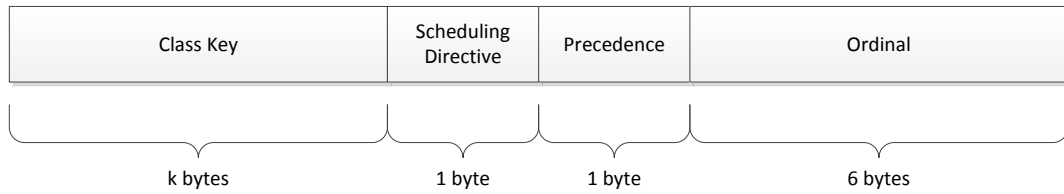


Figure 6. Example of Adaptive Heritrix frontier data structures and URL index

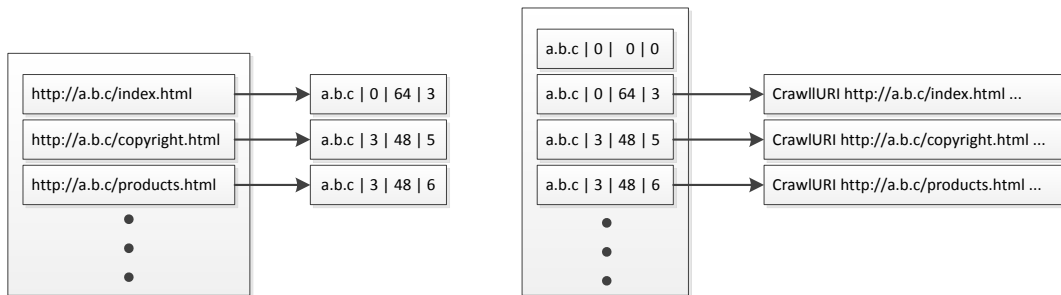
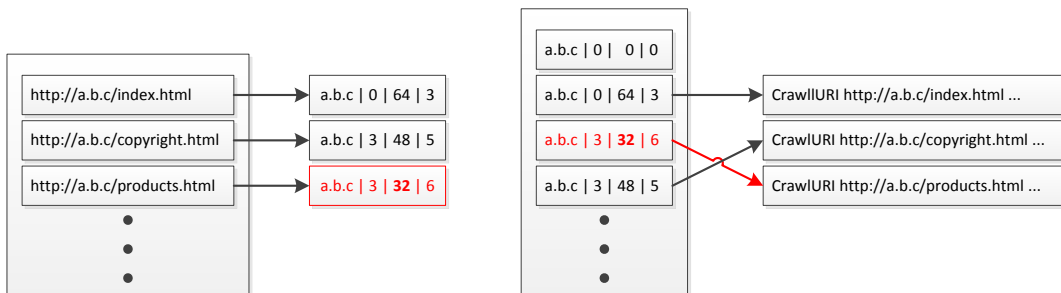


Figure 7. Example of updating the priority of a URL in the frontier of Adaptive Heritrix



393 When we need to update the priority of a URL already scheduled, we use this hash table to locate the
394 corresponding record from the frontier, update its priority, recalculate a key and insert it in the frontier
395 data structure at a new position. The fact that Heritrix employs an increasing counter to calculate the key
396 for each URL ensures that there are no collisions.

397 Figure 6 shows as an example the state of the frontier during a crawl for a domain `a.b.c`, where
398 there are three URLs queued for crawling. The first one is flagged for immediate downloading and has a
399 priority of 64, while the other two URLs have equal priority of 48. Upon an update of the priority of the
400 URL `http://a.b.c/copyright.html` from 48 to 32, the order of the two URLs is reversed, as
401 shown in Figure 7.

402 In the frontier we have developed, when a URL u is scheduled for crawling, first we have to check
403 whether the hash table mapping URLs to entries in the frontier contains an entry for u . If u is found, then
404 we update the priority, otherwise we need to check whether u has already been crawled.

405 The second feature we discuss enables Heritrix to receive prioritized URLs from other processes.
406 Heritrix provides an action directory, where processes having access to the same filesystem can write files
407 with seeds or URLs to be crawled. In order to fit Heritrix into the ARCOMEM crawling architecture and
408 receive URLs with priority scores from the online analysis phase, we have implemented a Web service
409 which receives scored URLs in an ARCOMEM-specific JSON format. In the simplest case, the required
410 information is an identifier for the crawl, the URL, a score in the range $[0, 1]$, and optionally a flag
411 indicating whether this URL should be blacklisted, i.e., not crawled at all. The developed Web service
412 enables Heritrix to receive links from any external process or even from other instances of Heritrix,
413 facilitating the distributed operation of the crawler. The URL score is transformed from the range $[0, 1]$
414 to an integer, as expected by Heritrix.

415 5. Evaluation

416 Since the ARCOMEM crawling architecture departs from the standard crawling architectures [7], it
417 is important to evaluate its impact on the effectiveness of a crawler. In this section we first evaluate
418 how adaptive and batch prioritization affects the performance of a crawler based on a set of simulation
419 experiments. Next, we compare a crawl performed by using the ARCOMEM crawling architecture to a
420 crawl performed by using the standard Heritrix crawler.

421 5.1. Adaptive and batch prioritization

422 We assume a baseline crawler implementing a best-first crawling strategy. We represent the topic of
423 the crawl with a topic vector, which is defined as follows. For each seed Web page, we download its
424 content and we create a term vector from it. The topic vector corresponds to the vector sum of the seed
425 page term vectors. We assume that URLs are prioritized according to their similarity to the topic vector.
426 More specifically, the priority of a URL u is computed as the average of: a) the cosine similarity between
427 the content of Web page p in which u was found and the topic vector and; b) the cosine similarity between
428 the anchor text of the out-link from p to u and the topic vector.

429 An adaptive crawler can update the score of an already scheduled Web page using a function such as
430 MAX, SUM, AVG. For example, the function MAX updates the priority of an already scheduled Web

431 page if the new priority was higher than the existing one. The function LAST always updates the score
 432 to the most recently computed one and the function FIRST is equivalent to the baseline crawler.

433 A crawler supporting batch prioritization schedules links for crawling only after having downloaded
 434 a batch of Web pages. In such a case, a URL can be discovered in many Web pages, so the cosine
 435 similarities are computed between the topic vector and the sum of the vectors of Web pages in which the
 436 URL was found, or the sum of the anchor text vectors associated with links pointing to the URL. In this
 437 setting, we also simulate a crawler that fetches the k pages with the highest priority from each domain,
 438 instead of fetching just one Web page with the highest priority.

439 To evaluate the focused crawler architectures, we perform simulated crawls on datasets created with
 440 three topics of DMOZ. We create three random samples of 20 seeds for each of the topics and the results
 441 we obtain for each configuration are the average of 9 simulations. For each set of seeds, we simulate
 442 a crawl of 10,000 Web pages. The topic vector we use to compute similarities between each topic and
 443 the crawled Web pages corresponds to the vector sum of the seed term vectors. For the evaluation of the
 444 results, we employ three measures: a) harvest ratio, which we define as the ratio of Web pages whose
 445 cosine similarity with the topic vector is greater than 0.333 over all crawled Web pages; b) average
 446 similarity of crawled pages; and 3) fraction of DMOZ subtopics with at least one crawled page.

447 Table 1 shows the evaluation results for adaptive prioritization with different priority update functions.
 448 The highest harvest ratio is achieved with the AVG function, while LAST achieves the highest fraction
 449 of DMOZ subtopics.

Table 1. Results of simulated adaptive crawls

Update function	Harvest ratio	Average Similarity	DMOZ topics
FIRST	0.3317	0.2945	0.4979
AVG	0.3609	0.3024	0.5779
MAX	0.3388	0.2967	0.5270
SUM	0.2679	0.2759	0.4650
LAST	0.3404	0.2961	0.5985

450 We also consider an additional parameter related to how the simulated crawler schedules links to
 451 crawl. For each queue, which corresponds to a domain, the crawler selects the bl URLs with the highest
 452 priority to crawl. In all previous experiments, $bl = 1$, meaning that each time the crawler selects the
 453 URL with the highest priority. For efficiency reasons, crawlers use higher values for the parameter bl . In
 454 the case of Heritrix, bl corresponds to the parameter *balance per queue*.

455 In our experiments, we test the values $bl = 1$ and $bl = 5$. The evaluation results for batch priority
 456 updating are shown in Table 2. We can observe that when increasing the value of bl , the effectiveness of
 457 the crawler drops, both in the case of a baseline crawler and an adaptive crawler (first 4 rows in Table 2).
 458 However, the combination of a higher bl value with batch updating performs better than an adaptive
 459 crawler with $bl = 5$ (rows 4 and 6 in Table 2).

460 5.2. Comparison of ARCOMEM versus standard crawl

Table 2. Results from simulated crawls with and without batch priority updating

Batch	<i>bl</i>	Update function	Harvest Ratio	Average Similarity	DMOZ topics
No	1	FIRST	0.3317	0.2945	0.4979
No	5	FIRST	0.2948	0.2819	0.4677
No	1	AVG	0.3609	0.3024	0.5779
No	5	AVG	0.3200	0.2897	0.5420
Yes	1	AVG	0.3556	0.3013	0.5260
Yes	5	AVG	0.3347	0.2952	0.5176

461 For evaluating the quality of the crawls in terms of how focused the collected documents are, we
 462 conducted a series of experiments. The main task of these experiments is to give an overview of the
 463 number of the crawled documents match the keywords given by the ICD, and how similar the textual
 464 content of the crawled documents is compared to the seed documents. Since there is no ground truth
 465 for the actual relevance of a document we have investigated some alternatives on how to measure the
 466 relevance of a document with respect to the given crawl definition. The dataset used for the experiment
 467 consists of two different crawls in the financial domain. One crawl was performed using the described
 468 ARCOMEM architecture while the other crawl is a standard Heritrix crawl.

469 The results we are presenting in this section are based on the document score generated by the Solr
 470 scoring module². This score is computed as follows. For each crawl, we create an inverted index of
 471 the crawled documents. Next, we form a query with the keywords provided by the ICD, or the most
 472 representative keywords from the textual content of the seed documents. The idea of these settings is to
 473 find documents with similar textual content to the ICD keywords or to the content of the seed documents.
 474 We assume that the documents matching the query are relevant and that their relevance is represented by
 475 their similarity to the query, or in other words, their retrieval status value (RSV).

476 Since the standard Solr scoring is based on TF/IDF, the whole document collection is taken into
 477 account for normalization. One drawback of this method may arise if all analyzed documents contain
 478 the relevant keywords. The IDF score which is used for normalization will be relatively low if many
 479 documents of the collection contain these words, and due to this, relevant documents may get a lower
 480 similarity score. In order to deal with this factor we checked how many of the documents contained
 481 our keywords. The results are shown in Table 3. The ARCOMEM crawl consist of 234,749 documents
 482 and the Heritrix Crawl contained 366,806 documents. In order to make both crawls comparable we
 483 analyzed the percentage of documents containing keywords from the ICD. The results show that most
 484 of the keywords do not appear inside most documents which allows us to use a TF/IDF based scoring
 485 approach. It becomes also visible that the percentage of occurrences of the keywords from the ICD in
 486 the crawled documents is in nearly all cases higher when crawling with the ARCOMEM framework than
 487 when crawling with Heritrix.

² <http://lucene.apache.org/core/>

Table 3. Percentage of crawled documents containing keywords from ICD

Keyword	ARCOMEM	Heritrix
Antonis Samaras	0.0175	0.0052
European Union	0.6011	0.4346
EU	1.7870	1.8986
financial crisis	0.8443	0.3645
Eurozone	0.3157	0.1238
bailouts	0.0826	0.0393
austerity	0.3838	0.1638
measures	1.4880	0.7993
strategy	1.9123	0.9362
growth	3.9319	1.9752
public	6.7357	3.5990
investments	1.3406	0.5545

488 The results of the experiments to assess the quality of the ARCOMEM and Heritrix crawls are shown
 489 in Tables 4 and 5, respectively. We compared four different settings for describing what makes a
 490 document relevant. These settings differ in how the terms are selected and the number of terms that
 491 are used for calculating the similarity score of crawled documents to the topic of the crawl. For the
 492 experiment in the first row (ICD), we use the terms from the ICD. For the experiments in the last three
 493 rows, we choose the terms based on their TF/IDF value for the seed documents; $maxQt=10$ indicates
 494 that the 10 most representative keywords are taken into account, which gives a very narrow definition
 495 of relevance. Additionally we used 50 and 100 words giving us a broader definition of relevance, since
 496 more words are considered to be important.

Table 4. Statistics on ARCOMEM Crawl

Ground Truth	Relevant Documents (%)	Average Similarity	Maximum Similarity	Standard Deviation
ICD	14.29	0.0178	0.8440	0.0426
Seeds $maxQt=10$	10.81	0.0179	0.4830	0.0263
Seeds $maxQt=50$	44.06	0.0049	0.4299	0.0149
Seeds $maxQt=100$	48.63	0.0060	0.5279	0.0172

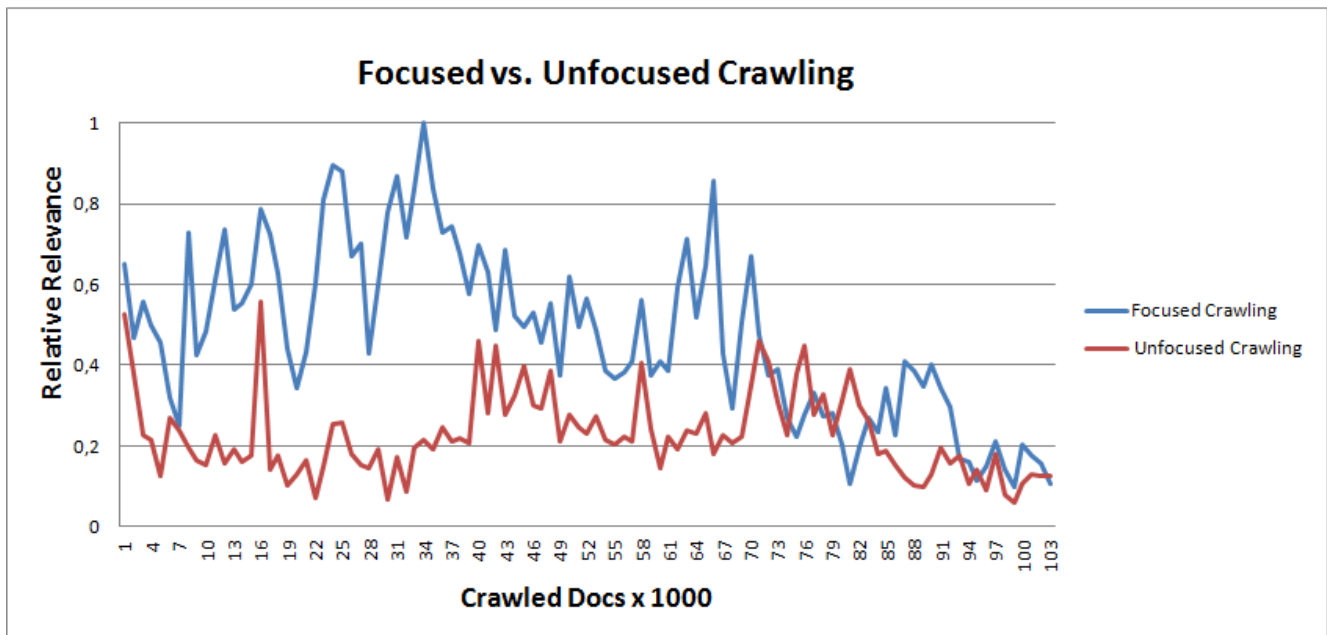
497 When comparing the results of the two crawls we see some obvious differences. The maximum
 498 similarity score per document was always the highest inside the ARCOMEM crawl, no matter which
 499 ground truth was chosen. For the first two setups, with a limited number of relevant keywords, the
 500 percentage of relevant documents was much higher inside the ARCOMEM crawl; when many different
 501 keywords are taken into account this changes. In contrast to that, the average similarity of the documents
 502 was higher for the Heritrix crawl when only few keywords are considered. This can be explained by

Table 5. Statistics on Heritrix Crawl

Ground Truth	Relevant Documents (%)	Average Similarity	Maximum Similarity	Standard Deviation
ICD	7.61	0.0196	0.7580	0.0440
Seeds maxQt=10	6.14	0.0175	0.2957	0.0264
Seeds maxQt=50	56.26	0.0026	0.3401	0.0114
Seeds maxQt=100	59.13	0.0021	0.3122	0.0099

503 looking at the absolute number of relevant documents for the different setups. The focused crawler
 504 found many more documents related to the small set of keywords. As a result, the average relevance has
 505 dropped, since this is calculated based only on the relevant documents.

506 Beside the results describing the relevance of all the crawled documents, we also analyzed how the
 507 relevance evolves over time. Figure 8 shows how the relevance of the crawled documents evolves over
 508 time. The relative relevance is calculated using the average similarity over 1000 crawled documents and
 509 dividing this value by the maximum of the averages.

Figure 8. Focused (ARCOMEM) and unfocused (Heritrix) crawling over time

510 We see that the relevance of the crawled documents exhibits large fluctuations during the crawl. While
 511 the Heritrix crawl does not show a certain tendency over the crawl, we see that the relative relevance of
 512 the documents of the ARCOMEM crawl increases up to the maximum after around 34,000 documents
 513 and then drops. Overall the relevance of the ARCOMEM based crawl was higher for most of the time.

514 Since one of the ideas of the ARCOMEM architecture is to use social networks to get additional
 515 content for the crawl we also used the described methods for analyzing the quality of links posted within
 516 Twitter. For this experiment we collected a total of 14,703 tweets related to our topic, out of which
 517 7,677 contained at least one URL. The content of these URLs was crawled and indexed in the same way

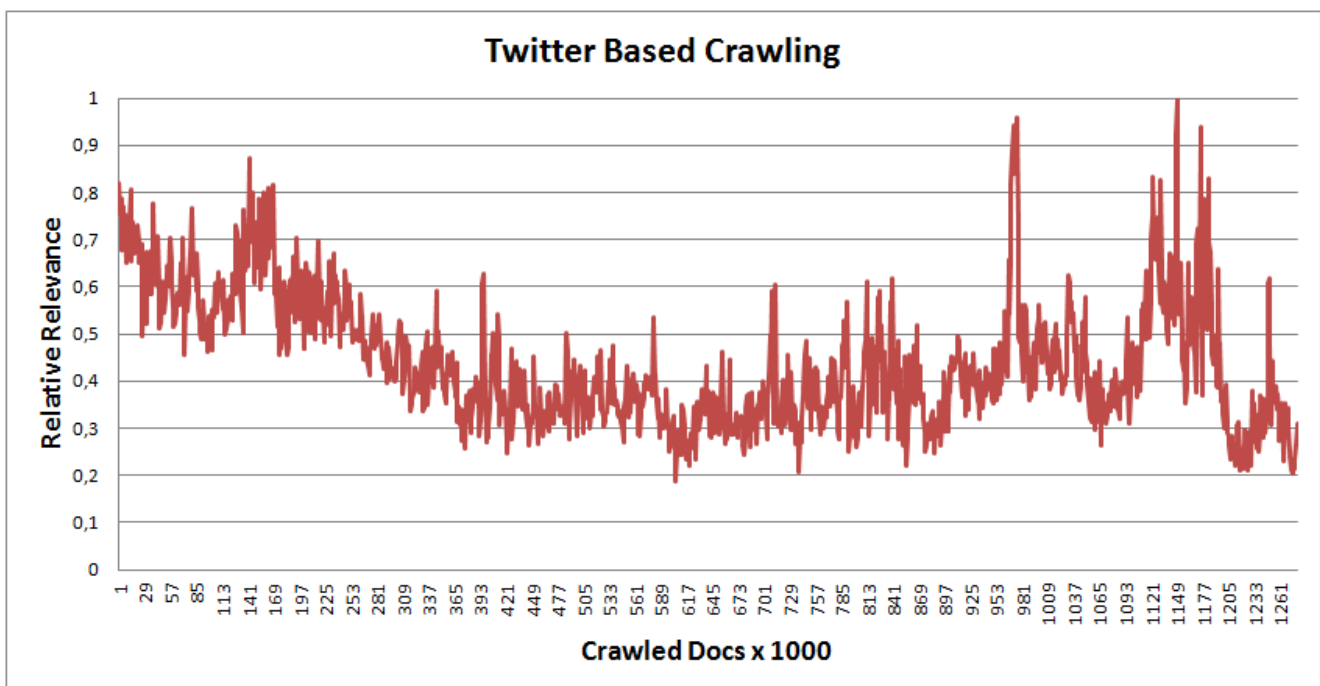
518 as we did it with the standard Heritrix and ARCOMEM crawl. Overall we performed this experiment
 519 on a set of 2.2 million crawled documents. The overall number of documents was much larger, but for
 520 this experiment we only took documents with textual content into account. Table 6 shows some basic
 521 statistics on the crawl.

Table 6. Statistics on Twitter based crawl

Ground Truth	Relevant Documents (%)	Average Similarity	Maximum Similarity	Standard Deviation
ICD	8.16	0.2435	2.1465	0.2299
Seeds maxQt=10	19.87	0.0348	1.7057	0.0805
Seeds maxQt=50	58.51	0.0180	1.5738	0.0489
Seeds maxQt=100	66.67	0.0193	1.6462	0.0493

522 We can see that the numbers for the percentage of relevant documents are comparable to the Heritrix
 523 and ARCOMEM crawls, except the number of relevant documents for the top 10 words from the seeds
 524 are much higher with 19%. The numbers for the average relevance and maximum relevance are not
 525 directly comparable with the results from the two other crawls because this crawl was stored in a different
 526 index and due to that the TF/IDF values are calculated based on different corpora. Nevertheless we can
 527 still see some parallels between the crawls: the average relevance is still the highest for documents
 528 related to the ICD and relative order of the other values is also comparable to the previous crawls.

Figure 9. Twitter based crawl over time



529 Figure 9 shows the evolution of the relative relevance of the crawled documents based on the 50 most
 530 representative terms from the seed set. Compared to the two previous crawls we cannot see the strong
 531 tendency of a dropping relevance over time.

532 6. Conclusions

533 The scale of the Web, the volatility of information found in it, as well as the emergence of social
534 media, require a shift in the way Web archiving is performed. Towards this goal, the ARCOMEM project
535 has developed a scalable and effective framework that allows archivists to leverage social media and
536 guide crawlers to collect both relevant and important information for preservation and future reference.

537 In this article, we have presented ARCOMEM's crawling architecture, providing a detailed
538 description of its main modules for extracting structured information from Web applications and
539 prioritizing the URLs to be crawled. We have also outlined the main features of a large-scale distributed
540 crawler, which can collect the content from billions of URLs while maintaining a high download rate.
541 The architecture we have described enables us to use either the large-scale crawler or an enhanced version
542 of Heritrix, for which we have described the required modifications we have implemented to support the
543 adaptive prioritization of URLs and the scheduling of URLs from remote processes.

544 Our experimental results show that the adaptive and batch prioritization, which are employed in the
545 proposed crawling architecture, are effective in acquiring relevant content. When comparing the quality
546 of a crawl performed with the ARCOMEM architecture against a crawl performed with Heritrix, we have
547 seen that the ARCOMEM crawler has downloaded earlier more relevant content. Overall, the proposed
548 crawling architecture is both extensible, by adding new modules to expand the analysis, and scalable,
549 offering a new approach to crawling content for Web archiving.

550 Acknowledgements

551 This work was funded by the European Commission under grant agreement n. 270239 (ARCOMEM).

552 Conflicts of Interest

553 Thomas Risse is co-editor of the special issue on Archiving Community Memories.

554 References

- 555 1. Koehler, W. A longitudinal study of Web pages continued: a consideration of document
556 persistence. *Inf. Res.* **2003**, *9*.
- 557 2. Twitter. Historical data not working. <https://dev.twitter.com/discussions/2483>, 2011.
- 558 3. Masanès, J. *Web Archiving*; Springer-Verlag New York, Inc.: Secaucus, NJ, USA, 2006.
- 559 4. Sigurðsson, K. Incremental crawling with Heritrix. Proceedings of the 5th International Web
560 Archiving Workshop (IWAW'05), 2005.
- 561 5. Risse, T.; Dietze, S.; Peters, W.; Doka, K.; Stavarakas, Y.; Senellart, P. Exploiting the Social and
562 Semantic Web for Guided Web Archiving. In *Theory and Practice of Digital Libraries*; Zaphiris,
563 P.; Buchanan, G.; Rasmussen, E.; Loizides, F., Eds.; Springer, 2012; Vol. 7489, pp. 426–432.
- 564 6. Plachouras, V.; Carpentier, F.; Masanés, J.; Risse, T.; Senellart, P.; Siehndel, P.; Stavarakas, Y. An
565 Architecture for Selective Web Harvesting: The Use Case of Heritrix. Proceedings of the 1st
566 International Workshop on Archiving Community Memories, 2013.
- 567 7. Olston, C.; Najork, M. Web Crawling. *Found. Trends Inf. Retr.* **2010**, *4*, 175–246.

- 568 8. Brin, S.; Page, L. The Anatomy of a Large-scale Hypertextual Web Search Engine. Proceedings
569 of the 7th International Conference on World Wide Web; Elsevier Science Publishers B. V.:
570 Amsterdam, The Netherlands, The Netherlands, 1998; WWW7, pp. 107–117.
- 571 9. Burner, M. Crawling towards eternity: Building an archive of the World Wide Web. *Web*
572 *Techniques Magazine* **1997**, *2*.
- 573 10. Heydon, A.; Najork, M. Mercator: A Scalable, Extensible Web Crawler. *World Wide Web* **1999**,
574 *2*, 219–229.
- 575 11. Najork, M.; Heydon, A. Handbook of Massive Data Sets; Kluwer Academic Publishers: Norwell,
576 MA, USA, 2002; chapter High-performance Web Crawling, pp. 25–45.
- 577 12. Shkapenyuk, V.; Suel, T. Design and implementation of a high-performance distributed Web
578 crawler. Proceedings of the 18th International Conference on Data Engineering, 2002, pp.
579 357–368.
- 580 13. Mohr, G.; Kimpton, M.; Stack, M.; Ranitovic, I. Introduction to heritrix, an archival quality web
581 crawler. Proceedings of the 4th International Web Archiving Workshop (IWAW'04); , 2004.
- 582 14. Boldi, P.; Codenotti, B.; Santini, M.; Vigna, S. UbiCrawler: A Scalable Fully Distributed Web
583 Crawler. *Softw. Pract. Exper.* **2004**, *34*, 711–726.
- 584 15. Lee, H.T.; Leonard, D.; Wang, X.; Loguinov, D. IRLbot: Scaling to 6 Billion Pages and Beyond.
585 *ACM Trans. Web* **2009**, *3*, 8:1–8:34.
- 586 16. Ntoulas, A.; Cho, J.; Olston, C. What's New on the Web?: The Evolution of the Web from a
587 Search Engine Perspective. Proceedings of the 13th International Conference on World Wide
588 Web; ACM: New York, NY, USA, 2004; WWW '04, pp. 1–12.
- 589 17. Fetterly, D.; Manasse, M.; Najork, M.; Wiener, J. A Large-scale Study of the Evolution of Web
590 Pages. Proceedings of the 12th International Conference on World Wide Web; ACM: New York,
591 NY, USA, 2003; WWW '03, pp. 669–678.
- 592 18. Cho, J.; Garcia-Molina, H. The Evolution of the Web and Implications for an Incremental
593 Crawler. Proceedings of the 26th International Conference on Very Large Data Bases; Morgan
594 Kaufmann Publishers Inc.: San Francisco, CA, USA, 2000; VLDB '00, pp. 200–209.
- 595 19. Olston, C.; Pandey, S. Recrawl Scheduling Based on Information Longevity. Proceedings of the
596 17th International Conference on World Wide Web; ACM: New York, NY, USA, 2008; WWW
597 '08, pp. 437–446.
- 598 20. Pandey, S.; Dhamdhare, K.; Olston, C. WIC: A General-purpose Algorithm for Monitoring Web
599 Information Sources. Proceedings of the 30th International Conference on Very Large Data
600 Bases. VLDB Endowment, 2004, VLDB '04, pp. 360–371.
- 601 21. Gouriten, G.; Maniu, S.; Senellart, P. Scalable, Generic, and Adaptive Systems for Focused
602 Crawling. Proceedings of 25th ACM Conference on Hypertext and Social Media; ACM: New
603 York, NY, USA, 2014; HT '13.
- 604 22. Tang, T.T.; Hawking, D.; Craswell, N.; Griffiths, K. Focused Crawling for Both Topical
605 Relevance and Quality of Medical Information. Proceedings of the 14th ACM International
606 Conference on Information and Knowledge Management; ACM: New York, NY, USA, 2005;
607 CIKM '05, pp. 147–154.

- 608 23. Menczer, F.; Pant, G.; Srinivasan, P.; Ruiz, M.E. Evaluating Topic-driven Web Crawlers.
609 Proceedings of the 24th Annual International ACM SIGIR Conference on Research and
610 Development in Information Retrieval; ACM: New York, NY, USA, 2001; SIGIR '01, pp.
611 241–249.
- 612 24. Chakrabarti, S.; van den Berg, M.; Dom, B. Focused crawling: a new approach to topic-specific
613 Web resource discovery. *Computer Networks* **1999**, *31*, 1623 – 1640.
- 614 25. Halkidi, M.; Nguyen, B.; Varlamis, I.; Vazirgiannis, M. THESUS: Organizing Web document
615 collections based on link semantics. *The VLDB Journal* **2003**, *12*, 320–332.
- 616 26. Ehrig, M.; Maedche, A. Ontology-focused Crawling of Web Documents. Proceedings of the
617 2003 ACM Symposium on Applied Computing; ACM: New York, NY, USA, 2003; SAC '03,
618 pp. 1174–1178.
- 619 27. Ahlers, D.; Boll, S. Adaptive Geospatially Focused Crawling. Proceedings of the 18th ACM
620 Conference on Information and Knowledge Management; ACM: New York, NY, USA, 2009;
621 CIKM '09, pp. 445–454.
- 622 28. Gao, W.; Lee, H.C.; Miao, Y. Geographically Focused Collaborative Crawling. Proceedings
623 of the 15th International Conference on World Wide Web; ACM: New York, NY, USA, 2006;
624 WWW '06, pp. 287–296.
- 625 29. De Bra, P.M.E.; Post, R.D.J. Information Retrieval in the World-Wide Web: Making Client-based
626 Searching Feasible. Selected Papers of the 1st Conference on World-Wide Web; Elsevier Science
627 Publishers B. V.: Amsterdam, The Netherlands, The Netherlands, 1994; pp. 183–192.
- 628 30. Hersovici, M.; Jacovi, M.; Maarek, Y.S.; Pelleg, D.; Shtalhaim, M.; Ur, S. The shark-search
629 algorithm. An application: tailored Web site mapping. *Computer Networks and ISDN Systems*
630 **1998**, *30*, 317–326.
- 631 31. Alpanidis, G.; Kotropoulos, C.; Pitas, I. Combining Text and Link Analysis for Focused
632 crawling-An Application for Vertical Search Engines. *Inf. Syst.* **2007**, *32*, 886–908.
- 633 32. Diligenti, M.; Coetzee, F.; Lawrence, S.; Giles, C.L.; Gori, M. Focused Crawling Using Context
634 Graphs. Proceedings of the 26th International Conference on Very Large Data Bases; Morgan
635 Kaufmann Publishers Inc.: San Francisco, CA, USA, 2000; VLDB '00, pp. 527–534.
- 636 33. Liu, H.; Janssen, J.; Milios, E. Using HMM to Learn User Browsing Patterns for Focused Web
637 Crawling. *Data Knowl. Eng.* **2006**, *59*, 270–291.
- 638 34. Partalas, I.; Paliouras, G.; Vlahavas, I. Reinforcement Learning with Classifier Selection for
639 Focused Crawling. Proceedings of the 2008 Conference on Artificial Intelligence; IOS Press:
640 Amsterdam, The Netherlands, The Netherlands, 2008; pp. 759–760.
- 641 35. Johnson, J.; Tsioutsoulouklis, K.; Giles, C.L. Evolving Strategies for Focused Web Crawling.
642 Proceedings of the 20th International Conference on Machine Learning; Fawcett, T.; Mishra, N.,
643 Eds. AAAI Press, 2003, pp. 298–305.
- 644 36. Zheng, H.T.; Kang, B.Y.; Kim, H.G. An ontology-based approach to learnable focused crawling.
645 *Information Sciences* **2008**, *178*, 4512 – 4522.
- 646 37. Bergman, M.K. *The Deep Web: Surfacing Hidden Value*, 2000.
- 647 38. Barbosa, L.; Freire, J. Siphoning Hidden-Web Data through Keyword-Based Interfaces.
648 Proceedings of the 19th Brazilian Symposium on Databases, 2004, pp. 309–321.

- 649 39. Madhavan, J.; Ko, D.; Kot, L.; Ganapathy, V.; Rasmussen, A.; Halevy, A. Google's Deep Web
650 crawl. *Proc. VLDB Endow.* **2008**, *1*, 1241–1252.
- 651 40. Arvidson, A.; Persson, K.; Mannerheim, J. The Kulturarw3 Project - The Royal Swedish Web
652 Archiw3e - An example of “complete” collection of web pages. Proceedings of the 66th IFLA
653 Council and General Conference, 2000.
- 654 41. Bailey, S.; Thompson, D. UKWAC: Building the UK's First Public Web Archive. *D-Lib*
655 *Magazine* **2006**, *12*.
- 656 42. Cathro, W.; Webb, C.; Whiting, J. Archiving the Web: The PANDORA Archive at the National
657 Library Australia. *National Library of Australia Staff Papers* **2009**.
- 658 43. Spaniol, M.; Denev, D.; Mazeika, A.; Weikum, G.; Senellart, P. Data Quality in Web Archiving.
659 Proceedings of the 3rd Workshop on Information Credibility on the Web; ACM: New York, NY,
660 USA, 2009; WICOW '09, pp. 19–26.
- 661 44. Denev, D.; Mazeika, A.; Spaniol, M.; Weikum, G. SHARC: Framework for Quality-conscious
662 Web Archiving. *Proc. VLDB Endow.* **2009**, *2*, 586–597.
- 663 45. Gomes, D.; Miranda, J.a.; Costa, M. A Survey on Web Archiving Initiatives. Proceedings
664 of the 15th International Conference on Theory and Practice of Digital Libraries: Research
665 and Advanced Technology for Digital Libraries; Springer-Verlag: Berlin, Heidelberg, 2011;
666 TPDFL'11, pp. 408–420.
- 667 46. Tahmasebi, N.; Gossen, G.; Kanhabua, N.; Holzmann, H.; Risse, T. NEER: An Unsupervised
668 Method for Named Entity Evolution Recognition. Proceedings of the 24th International
669 Conference on Computational Linguistics; Kay, M.; Boitet, C., Eds. Indian Institute of
670 Technology Bombay, 2012, COLING' 12, pp. 2553–2568.
- 671 47. Plachouras, V.; Stavarakas, Y.; Andreou, A. Assessing the Coverage of Data Collection Campaigns
672 on Twitter: A Case Study. OTM Workshops; Demey, Y.T.; Panetto, H., Eds. Springer, 2013, Vol.
673 8186, *Lecture Notes in Computer Science*, pp. 598–607.
- 674 48. Faheem, M.; Senellart, P. Intelligent and Adaptive Crawling of Web Applications for Web
675 Archiving. Proceedings of the 13th International Conference on Web Engineering (ICWE);
676 Daniel, F.; Dolog, P.; Li, Q., Eds. Springer Berlin Heidelberg, 2013, Vol. 7977, *Lecture Notes in*
677 *Computer Science*, pp. 306–322.
- 678 49. Faheem, M.; Senellart, P. Demonstrating intelligent crawling and archiving of web applications.
679 Proceedings of the 22nd ACM International Conference Information and Knowledge Manage-
680 ment; ACM: New York, NY, USA, 2013; CIKM '13, pp. 2481–2484.
- 681 50. Cunningham, H.; Maynard, D.; Bontcheva, K.; Tablan, V.; Aswani, N.; Roberts, I.; Gorrell, G.;
682 Funk, A.; Roberts, A.; Damljanovic, D.; Heitz, T.; Greenwood, M.A.; Saggion, H.; Petrak, J.; Li,
683 Y.; Peters, W. *Text Processing with GATE (Version 6)*; 2011.