

Characterizing a National Community Web

Daniel Gomes

and

Mário J. Silva

Universidade de Lisboa, Faculdade de Ciências

Portugal

This paper presents a characterization of the community web of the people of Portugal. We defined criteria for delimiting this web, based on our past experience of crawling pages related to Portugal, and collected over 3.2 million documents from 46,000 sites satisfying those criteria. Our characterization was derived from this crawl. We describe the rules that we established for defining the boundaries of this community web and the methodology used to gather statistics. Statistics cover: the number and domain distribution of sites; the number, type and size distribution of text documents; and the linkage structure of this web. We also show how crawling constraints and abnormal situations on the web can influence the statistics.

Categories and Subject Descriptors: H.3.5 [**Information Storage and Retrieval**]: Online Information Services—*Web-based services*; H.3.7 [**Information Storage and Retrieval**]: Digital Libraries—*Collection*; C.2.5 [**Computer-communication Networks**]: Local and Wide-Area Networks—*Internet*

General Terms: Statistics, Web, Internet, Portugal

Additional Key Words and Phrases: Web Characterization, Portuguese Web, Web Measurements, Web Communities

1. INTRODUCTION

A characterization of the web is of great importance. It reflects technological and sociological aspects and permits to understand how the web has evolved. An accurate characterization of the web enables improving the design and performance of applications that use the web as a source of information (e.g. crawlers, proxies, search engines) [Cho and Garcia-Molina 2000].

The web can be characterized from multiple perspectives using numerous metrics. This is a challenging task, mainly because of its large dimension and permanent evolution [Leung et al. 2001]. Producing a feasible general characterization is hard, and some statistics derived from the analysis of the global web may not hold as we scale down to more restricted domains. The web has partitions with specific characteristics that, given their small presence, do not become visible in a general

Authors' address: Faculdade de Ciências da Universidade de Lisboa, Departamento de Informática, Campo Grande, 1749-016 Lisboa, Portugal;
email: dcg,mjs@di.fc.ul.pt.

Permission to make digital/hard copy of all or part of this material without fee for personal or classroom use provided that the copies are not made or distributed for profit or commercial advantage, the ACM copyright/server notice, the title of the publication, and its date appear, and notice is given that copying is by permission of the ACM, Inc. To copy otherwise, to republish, to post on servers, or to redistribute to lists requires prior specific permission and/or a fee.

© 2004 ACM 0000-0000/2004/0000-00000 \$5.00

ACM Journal Name, Vol. 1, No. 1, February 2004, Pages 0–0??.

web characterization. However, these partitions can be of interest to relatively large communities, such as those representing national or cultural groups. Additionally, characterizing a small partition of the web is quite accessible and can be done with great accuracy.

In this paper, we present a detailed characterization of a national community web. This work was conducted in the context of a study of the Portuguese web, broadly defined as the set of pages of cultural and sociological interest to the people of Portugal. The results are derived from a crawl performed by tumba!, a search and archival engine for the Portuguese web [Silva 2003]. We focused our study on textual contents available on the Portuguese web, identifying metrics that would help us in the design and improvement of our system. The statistics themselves are interesting to anyone who manipulates these data or will compare it with our snapshot in the future. We compare our results with related work. However, we need to be cautious about the conclusions drawn, because the results were gathered during different periods and using distinct methodologies, which often are not detail enough. The identification of the meaningful statistics for a community web characterization and the methods used to gather and interpret the collected data could be useful to a wider audience. We detail our crawling policy and show how the crawling and data analysis processes can strongly influence the obtained statistics.

This paper is organized as follows: the remaining of this section presents the adopted terminology. Section 2, presents our heuristics for defining the boundaries of the Portuguese web. In the following 2 sections, we present the crawler configuration and the crawling results. In sections 5, 6 and 7 we describe the statistics related to web sites, documents and structure, respectively, derived from the crawl. Section 8 introduces related work. Finally, in section 9 we draw our conclusions and present directions for future research.

1.1 Terminology

The concepts used in this study were adapted from the terminology proposed by the W3C [W3C 1999].

- Publisher: entity responsible for publishing information on the web;
- Document: file resulting from a successful HTTP download;
- Page: HTML document;
- Web Site: collection of documents referenced by URLs that share the same host name, (a discussion about the definition of web site can be found in [O’Neill 1999]).
- Host page: document identified by an URL where the file path component is empty or a ‘/’ only.
- Subsite: cluster of documents within a web site, maintained by a different publisher than that of the parent site.
- Host aliases: sites that have different names but are hosted on the same IP address and have the same host page.

2. IDENTIFYING THE BOUNDARIES OF A COMMUNITY WEB

The web is designed to break all the geographical barriers and make information universally available. However, as the web is the product of multiple user groups, it is possible to identify partitions within it containing the sites of interest to these groups. These are designated as community Webs and can be defined as the set of documents that refer to a certain subject or are of interest to a community of users.

Detection of a community web is not always obvious, despite of various existing methods that can be used to identify its sites.

If we are interested in a small and static set of documents, then enumerating all the documents that compose the community web can be adequate. However, it becomes very expensive to maintain the list of documents if it grows or changes frequently [Webb 2000].

We can also use the link structure [Flake et al. 2000] of the web, but we'll find difficulties identifying documents loosely interlinked, even if they refer to the same subject. For instance, the sites of several concurrent companies in the same business, will not likely link to each other.

We can identify documents related to a country through a common country code Top Level Domain (ccTLD) [Postel 1994; Zabicka 2003]. In this case, we would exclude all the documents related to that country hosted under a domain outside the ccTLD. On the other hand, this rule would also include sites not related to the country, but being hosted under its ccTLD. For instance, multi-national companies commonly register their name under many domains to protect their brands.

The language in which the documents are written is a good indicator of which country they are related [Albertsen 2003]. However, problems arise if the language is not exclusive to a single country: we couldn't include all the documents written in English within a British community web.

As a result, a precise definition of which documents should constitute a community web is in general hard to obtain and is conditioned by the rules and resources used.

The community web of our study is the Portuguese web. We define it as the set of documents containing information related to Portugal or of major interest to the Portuguese people. Our first approach for establishing the boundary of the Portuguese web outside the .PT domain was to harvest all the documents written in the Portuguese language. Soon, we found that this would impose the downloading of a large number of documents, specially from Brazil, containing information not highly related to Portugal. As defining rule, we consider as part of the Portuguese web those documents that satisfy one of the following conditions:

- 1. Hosted on a site under the .PT domain;
- 2. Hosted on a site under the .COM, .NET, .ORG or .TV domains, written in the Portuguese language and with at least one incoming link originating in a web page hosted under the .PT domain.

This definition aims to be easily set as a crawling policy and guarantees that the crawler does the best coverage of the Portuguese web.

Condition 1 intends to include the sites that constitute the core of the Portuguese web. A list of the most popular sites, accessed from the homes of a panel of

Portuguese users, during 2002 and 2003 [Marktest 2003] showed that 49.5% of the sites were hosted under the .PT domain. Based on this information we considered the sites hosted under the .PT domain as the core of the Portuguese community web.

Condition 2 intends to include the increasing number of Portuguese sites that are registered outside the .PT domain [Zook 2000]. Previous work [Flake et al. 2000; Gibson et al. 1998] showed that the link structure of the web can be used to define communities. We assumed that the probability of a site hosted outside the .PT domain belonging to the Portuguese web community decreases as the number of hops in the web graph to the core increases. So, in order to restrict the inclusion of sites outside the .PT domain to the ones with the highest probability of being part of the Portuguese web, we limited the number of hops to 1. Condition 2 imposes that only documents directly linked from a site hosted under the .PT domain are part of the Portuguese community web. However, Brazilian documents linked from the .PT domain and hosted under the allowed domains, such as .COM will still be considered as part of the Portuguese web.

The definition of the Portuguese community web has an implicit geographical context. We ran an experiment with the purpose of comparing the coverage of the Portuguese community web outside the .PT domain (condition 2) by different alternatives.

We gathered a list of 25 Portuguese sites hosted outside the .PT domain suggested by Portuguese users. Then we examined the suggested sites and verified that they were part of directly related to the Portuguese web community. All the sites were written in Portuguese and referred to distinct subjects such as sports, humour or radio. We compared our proposed defining criterion against 2 alternatives, based on tools that provide geographic context data for web sites.

In the first alternative we used 2 commercial tools, Ip2location [Center 2003] and Maxmind [LLC 2003], and extracted a geographical location for each site on the list. We considered a site as part of the Portuguese web if the tool returned that the site was located in Portugal. For 2 submissions of the same site, Maxmind returned different results. Except for this situation, both tools presented the same results, which lead us to believe that they are based on the same data.

Our second alternative was to access a whois database [Harrenstien et al. 1985] to identify the Portuguese sites hosted outside the .PT domain. For each site we obtained the contact address of the correspondent domain registrant. If this address was located in Portugal, we considered the site as part of the Portuguese web.

Finally, following our proposed definition, we checked if the sites had at least one link from a site hosted under the .PT domain. We used the search engines Google [Google 2003] and AllTheWeb [Overture Services 2003], to identify pages that link to the sites.

We also tried to obtain geographic information through the DNS LOC record [Davis et al. 1996] but none of domains had an associated record of this kind.

Table I presents the results obtained through the three definitions of the Portuguese web. The geographical tools identified only 44% of the Portuguese sites, although they always returned an answer to the location requests. 76% of the Portuguese sites were identified through the registrant information. For 24%, of the

Definition	% sites identified	% information unavailable
Geographical tools	44	0
Whois registrant address	76	24
Linked from .PT	82	12

Table I. Comparison between alternative definitions of the Portuguese web.

sites the whois database didn't contain the information regarding the requested domain. For some of these cases, we found the registrant information on another whois server. As there isn't a central whois database, the registrants information is distributed over the several registrars, which causes inconsistencies among whois databases. The registrant address revealed to be a precise method to identify Portuguese sites outside the .PT domain. All the sites in our list which had a whois record available were correctly identified. Therefore, the whois databases could be the solution to the problem of distinguish Brazilian sites from Portuguese sites outside the .PT domain. However, most of the whois databases are not publicly available or explicitly forbid their access by automated programs, which collides with our purpose of having a definition of the Portuguese web that can be implemented as a crawling policy. The existence of several record formats also makes it difficult to automatically process whois records. Additionally, companies that provide hosting services support several distinct sites identified by sub-domains and the whois registries only keep information about second-level domains. This does not include the sub-domains of hosted web sites. This is a serious restriction if we consider, for instance, all the Portuguese blogs hosted under `blogspot.com`. If we had followed this approach, many Portuguese web sites hosted outside Portugal would not be considered as part of the Portuguese web.

We observed that 82% of the suggested sites would be included in the Portuguese web using our criteria: they were written in the Portuguese language and had at least one link from a site hosted under a .PT domain. The results show that our proposed definition of the Portuguese web provides the best coverage of the suggested sites.

3. CRAWLER CONFIGURATION

A crawler begins its task of harvesting the web collecting the documents referenced by a initial set of URLs, called the seeds. Then it iteratively extracts links to new URLs and collects their contents.

Crawlers are configured or developed according to the purpose of the data they gather. A crawler of a large scale search engine aims to collect pages with the highest PageRank [Cho et al. 1998; Brin and Page 1998]. On the other hand, archive crawlers focus on crawling the most pages on a given partition [Day 2003]. In our study, we configured Viúva Negra (VN) [Gomes 2003], the web crawler of the tumba! search engine, to get the most information possible about the Portuguese web. We initialized it with a set of 112,146 seeds gathered from previous crawls and user registrations that include all the hosts registered under the .PT domain. We imposed on it the minimum constraints that ensure an acceptable performance of the crawler, considering the resources available and the need to make it robust against the usual anomalies in the traversed web graph [Henzinger 2003]. A docu-

ment was considered to be valid if it was part of the Portuguese web, as defined in the previous section. In addition, the following crawler conditions had to be met:

- multiple text types: we considered not only documents of the text MIME type but also documents of common MIME application types that we could convert to text. Accepted MIME types are: text/html, text/richtext, text/tab-separated-values, text/plain, text/rtf, application/pdf, application/rtf, application/x-shockwave-flash, application/x-tex, application/msword, application/vnd.ms-excel, application/excel, application/mspowerpoint, application/powerpoint and application/vnd.ms-powerpoint.
- URL depths less than 6: the crawler followed at most 5 links in breadth-first search order, from the seed of the site until it reached the referenced document. When crawling a site, we considered that any link found to a different site would be set as a seed to that site. This way, we guaranteed that any page with a link originated on a .PT domain would be visited, including Portuguese subsites hosted on foreign sites. Consider for instance, the site www.yahoo.com and its Portuguese subsite www.yahoo.com/users/myPortugueseSite/. If the crawler had visited only the seed www.yahoo.com, it would not have identified that the site wasn't part of the Portuguese web, and exited without finding the Portuguese subsite;
- documents downloaded in less than 1 minute: this prevents very slow web servers from blocking the progress of the crawl;
- document size under 2 MB: this prevents the download of very large files available on the web, such as database dumps.

3.1 Avoiding traps

A crawler trap is a set of URLs that cause a crawler to traverse a site indefinitely. They are easily noticed due to the large number of documents discovered in the site [Heydon and Najork 1999]. In order to prevent the crawling of infinite sites, we set VN to visit a maximum of 8000 URLs per site. This turned out to be an acceptable limit, considering the dimensions of the Portuguese web sites (see section 5). This constraint reduced the number of unnecessary downloads and increased the robustness of the crawler, but it wasn't enough to prevent traps from biasing a web characterization. We found that most of the traps are unintentional, being caused mainly by session identifiers embedded in the URLs, or poorly designed HTTP web applications that dynamically generate an infinite number of URLs that reference a small set of documents.

This raises the issue of how should these documents be considered in a characterization. They should not be excluded because they are available online and represent part of the web. However, we can not let them bias a characterization due to its "infinite" presence. We adopted the solution of setting VN as a very patient web surfer as a compromise. After seeing the same bitwise identical document 50 times, VN gives up on following links for that site, keeping all the information crawled until then. This limitation intends to avoid spider traps that always return the exact same content. If the trap generates slightly different contents, we identify it when the site reaches the maximum number of documents allowed. A criterion that identifies documents with distinct URLs and contents as being similar enough to be

State	# URLs	%
200	3235140	83.9
302	193870	5.0
404	132834	3.4
TimedOut (-8)	45486	1.2
301	39920	1.0
ExcludedByREP (-2)	35596	0.9
500	33247	0.9
NotAllowedType (-5)	25976	0.7
403	18598	0.5
UnknownHost (-14)	17842	0.5
SizeTooBig (-4)	17453	0.5
ConversionError (-11)	13986	0.4
Other	23244	0.6
Total	3856436	100.0

Table II. Summary of the status codes associated to the URLs visited. The positive numbers represent the HTTP response codes and the negative numbers represent VN special codes that identify the reason why the contents referenced by the URLs were not collected.

considered the same is highly subjective. If several documents are similar except in a banner ad that changes on every download, they could reasonably be considered the same. However, when the difference between them is only as short as the licitation value on an online auction, the small difference could be very significant! As a result, we considered the computation of partial similarity between documents to be too expensive and risky to be applied in the identification of spider traps during the crawling process.

We didn't identified any intentional trap in our crawl. However, they can be created using DNS wildcarding [Barr 1996] to resolve any possible host name within a domain to the same IP address, generating an infinite number of host aliases and giving web crawlers the illusion that each site serves only a small number of pages. In order to mitigate this situation, VN avoids to crawl host aliases by identifying them through a pre-computed list gathered from a previous crawl. We verified the host aliases list using online information (IP address and host pages) before starting the crawl, so that it could be as accurate as possible. Hence, there are host aliases that have disappeared since the previous crawl.

4. CRAWLING STATISTICS

The results presented were extracted from a crawl performed between the 1st of April and the 15th of May, 2003. VN visited a total of 146076 sites, processed over 3.8 million URLs and downloaded 78 GB of data¹.

Table II presents the statistics of the download status of crawled URLs. Almost 84% of the requests resulted in a successful download and only 3.4% resulted in a 404 (File Not Found) response code, which indicates that most of our seeds were valid and that broken links are not as frequent in this web as reported in other studies

¹The information gathered in this crawl is available for research purposes at http://xldb.fc.ul.pt/linguateca/WPT_03.html.

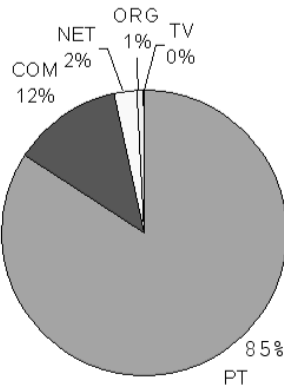


Fig. 1. Distribution of sites per Top Level Domain.

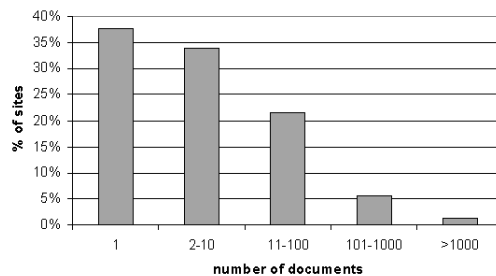


Fig. 2. Distribution of documents per site.

[Najork and Heydon 2001; Spinellis 2003]. There were over 6% of redirections and the crawler failed to fetch and parse a document within 1 minute in 1.2% of the requests. The Robots Exclusion Protocol prevented VN from downloading 0.9% of the URLs, and about the same number of URLs resulted in an Internal Server Error (500). The number of documents with a not allowed MIME type (0.7%) is underestimated, because extracted links that had names hinting that the referenced content didn't belong to one of the allowed types (ex. files with a .JPEG extension) were not crawled. The UnknownHost error (0.5%) is caused by URLs referencing host names that no longer have an IP associated. We found that only 0.5% of the referenced files had a size bigger than 2 MB and the conversion to text was not possible in 0.4% of the cases. The remaining situations (0.6%) included other HTTP response codes, unidentified errors, socket and connection errors; each of these represents less than 0.1% of the total number of downloaded documents.

5. SITE STATISTICS

We considered that a site is part of the Portuguese web if it hosted at least one document considered as part of the Portuguese web. We identified 46457 sites as being part of the Portuguese web. 85% of the sites were under the .PT domain,

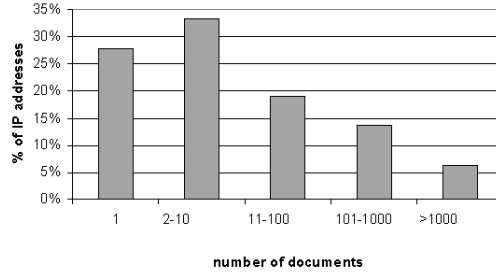


Fig. 3. Distribution of documents per IP address.

# sites per IP	# IP addresses	%
1	4643	67.7
2-10	1931	28.2
11-100	247	3.6
101-1000	30	0.4
>1000	5	0.1
Total	6856	100.0

Table III. Distribution of the number of sites hosted per IP address.

12% were under the .COM, 1% were under the .ORG domain and just 3 sites were under the .TV (see Figure 1). 60% of the web sites names started with "WWW". A Portuguese web site has an average of 70 documents, but the size distribution is very skewed, as shown in Figure 2. We were surprised by the high number (38%) of sites having a single document. We visited a random sample of these sites and observed that most warned readers that they were under construction, or that the site moved to a different location. We also found a few cases where the host page was completely written using scripting languages from which our parser couldn't extract links. A typical web site had less than 101 documents (93%); 6% had between 101 and 1000 documents and only 1% of the sites had more than 1000 documents. We identified only 577 sites that hosted more than the maximum number of 8000 documents. We observed that most of these sites were huge database dumps available online through dynamically generated web pages. We concluded that despite the restriction on the site size, we were able to exhaustively crawl 99% of the sites.

The distribution of documents per IP address is more uniform (see Figure 3). The percentage of IP addresses that host just one document is 28%, IP addresses that host 2 to 10 documents represent 33%, those which host between 11 and 1000 documents represent 33% and only 6% host more than 1000 documents.

Table III shows that over 32% of the IP addresses host more than 1 site. Each IP address hosts an average of 6.78 sites. Silva et al. [Silva et al. 2002b] compared results from 2 crawls of the .PT domain performed in 2001 and 2002, and observed that the number of sites per IP address grew from an average of 3.78 to 4.57 sites per IP address. Our result suggests that this number continues to grow. There are

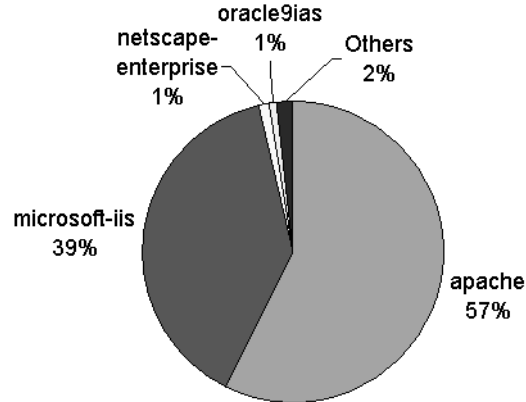


Fig. 4. Distribution of web servers.

5 IP addresses that host more than 1000 sites. These 5 IP addresses are from web portals that offer their clients a virtual host under the portal domain, providing a proper host name for their site, instead of having it as a subsite. Virtual hosts are very popular on the Portuguese web: 82% of all sites are virtual hosts. It is important to distinguish host aliases from distinct virtual hosts. The first occur when multiple names refer to the same site, for instance `http://xldb.fc.ul.pt` and `http://xldb.di.fc.ul.pt`. Distinct virtual hosts are distinct sites hosted on the same machine, such as `http://xldb.di.fc.ul.pt` and `http://lasige.di.fc.ul.pt`. In our crawl we found out that 8.5% of the virtual hosts were host aliases.

5.1 Web Servers

We identified 172 distinct HTTP web servers, figure 4 presents their distribution. The Portuguese sites are mainly hosted at Apache (57%) and Microsoft IIS (39%) web servers. The next two web servers (netscape-enterprise and oracle9ias) represent just 1% each and the remaining just 2%. Statistics on the global web present a similar percentage of Apache web servers (62.57%), but a considerably smaller percentage of Microsoft IIS servers (27.45%) [Netcraft Ltd 2004]. On the other hand our distribution of web server software contrasts with the one obtained by Boldi et al. [Boldi et al. 2002] for the Africa web, in which there is a dominance of Microsoft IIS over Apache (56.1% against 37.95%).

Security experts encourage webmasters to don't provide the Server HTTP field or to provide wrong answers in order to mislead possible attackers. From our experience we believe that these recommendations are usually not followed by the Portuguese webmasters. However, if they become popular, it will be very difficult to correctly identify the distribution of web server software.

6. DOCUMENT STATISTICS

In this section, we present metrics regarding the length of URLs, MIME types, size, language and meta-data of documents.

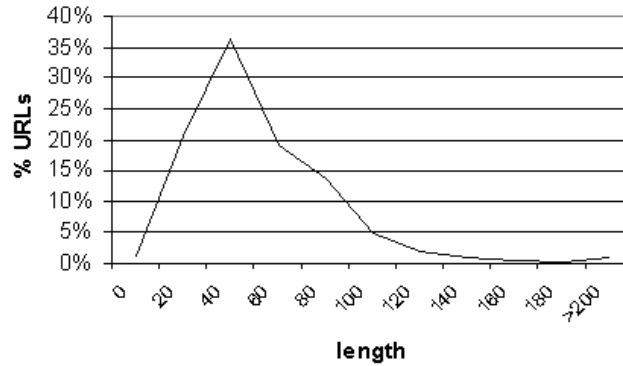


Fig. 5. Distribution of URL lengths.

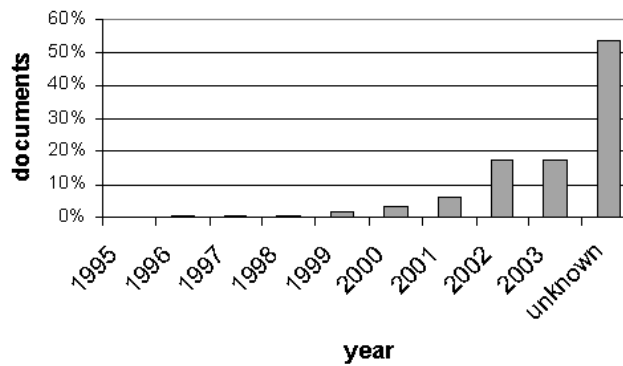


Fig. 6. Distribution of Last-Modified Dates.

6.1 URLs

Every web application must have some kind of data structure that maps into URLs. However, we didn't find in the literature a study discussing the lengths of the URLs. Nowadays, the size of URLs is in practice unlimited. We found valid URLs with lengths varying from 5 to 1368 characters. Figure 5 shows the distribution of URL lengths (not considering the initial 7 characters of the protocol) over the number of the documents. Most of the documents have an URL length between 20 and 100 characters, with an average value of 62 and median of 54. Analyzing the URLs we found that 2.3% contained parameters suggesting that the referenced document had been dynamically generated.

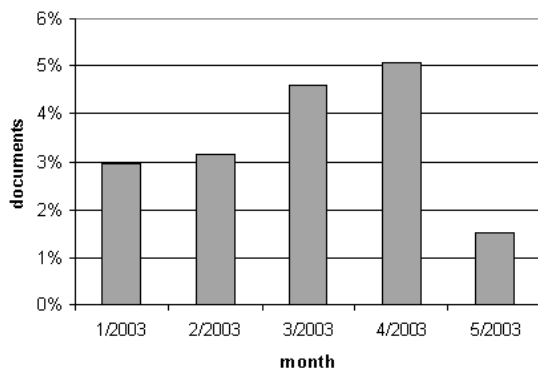


Fig. 7. Distribution of Last-Modified Dates in the last 4 months.

6.2 Last Modified Date

HTTP provides a header field (Last-Modified Date) that should indicate the date of the last modification of documents. However, as shown in Figure 6, most of the documents (53.5%) returned an unknown value for this field. Plus, Mogul showed that even the returned values are many times inaccurate due to the incorrectly set web server clocks (among other situations) [Mogul 1999b]. An analysis of the URLs with unknown values revealed that 82% of them had embedded parameters. We speculate that most of them are recent and would significantly increase the percentage of documents modified in the last months (see Figure 7), since mechanisms to dynamically generate documents are usually used to reference short life contents, such as news.

We believe that the last-modified header is a weak metric for evaluating changes and evolution of contents on the web, so metrics like these are meaningful only in the context of analysis of consecutive crawls [Wills and Mikhailov 1999; Fetterly et al. 2003].

6.3 MIMEs & Sizes

The rightmost column of Table IV shows the distribution of documents per MIME type, (we grouped all the MIME types used for Microsoft Powerpoint files under the name *powerpoint* and all the Microsoft Excel files under the name *excel*). The predominant text format is text/html, present in over 95% of the collected documents, followed by application/pdf with just 1.9%.

In our first approach to determine the size of the documents, we analyzed the values of the HTTP header field Content-Length but we noticed that 33% of the documents returned an unknown value. We then recomputed our results replacing the unknown sizes by the sizes of the documents. The differences on the average sizes between the results were insignificant except for text/html where the size grew from 12.2 KB to 20.5 KB. In Table V, the second and third columns show the average sizes of documents and corresponding extracted texts (without any formatting tags), and the fourth column presents the ratio between the length of

MIME	# documents	%
text/html	3104771	95.9
application/pdf	62141	1.9
text/plain	33091	1.0
application/x-shockwave-flash	17598	0.5
application/msword	14014	0.4
powerpoint	2085	0.1
excel	915	0.0
application/x-tex	222	0.0
text/rtf	194	0.0
application/rtf	66	0.0
text/tab-separated-values	41	0.0
text/richtext	2	0.0
Total	3235140	100.0

Table IV. Number of documents and relative presence on the web for each MIME type collected.

MIME	avg Doc Size(KB)	avg Text Size(KB)	% text
powerpoint	1054.9	7.0	1
text/rtf	475.6	1.2	0
application/pdf	207.4	13.6	7
application/rtf	121.3	4.7	4
application/msword	118.6	9.9	8
excel	50.4	21.9	43
application/x-shockwave-flash	43.9	0.3	1
text/html	20.5	2.5	12
text/richtext	16.3	16.2	99
application/x-tex	16.1	14.7	91
text/plain	10.5	7.8	74
text/tab-separated-values	3.9	3.8	97

Table V. Average size, extracted text size, percentage of extracted text.

the extracted text and document size. We can see that the size of the documents is almost inversely proportional to the size of the texts extracted. A curious fact is how documents of text/plain result in just 74% of text. We analyzed some of these documents and discovered that some web servers return text/plain, when the file type of the document is not recognized. Therefore, some PowerPoint Presentation files (.PPS) or Java Archives (.JAR) were incorrectly processed as text/plain, resulting in poor extraction of text from these files.

Figure 8 shows the general distribution of document sizes. Most documents have between 4 and 64 KB. The mean size of a document is 32.4 KB and the mean size of the extracted texts is 2.8 KB. The total size of the documents was 78.430 GB, while the total size of extracted texts was just 8.791 GB.

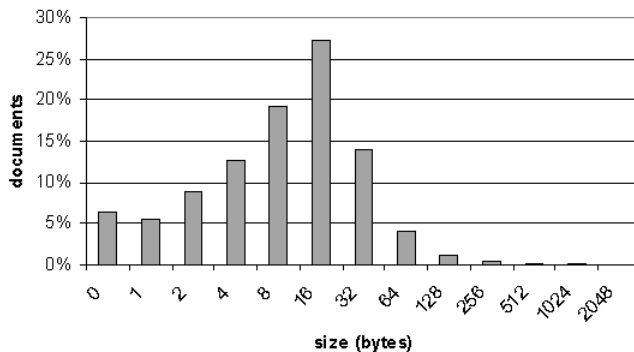


Fig. 8. Distribution of document sizes.

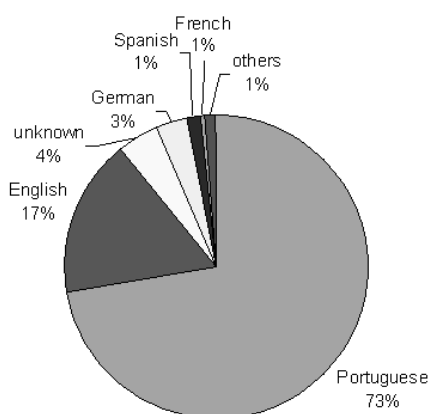


Fig. 9. Distribution of Languages.

6.4 Language Distribution

Our crawler can identify the language of collected documents based on an idiom detector that implements an n-gram algorithm [Cavnar and Trenkle 1994]. Figure 9 shows the distribution of languages on the documents of the Portuguese web (including documents written in all languages hosted under the .PT domain): 73% of the documents were written in Portuguese, 17% in English, 3% in German, 1% in Spanish, 1% in French and 1% in other languages. According to O’Neill et al., on the global web 72% of the pages are written in English and only 2% are written in Portuguese [O’Neill et al. 2003].

Identifying the language of a document is sometimes a hard task because there are documents with short text or written in several languages. The idiom detector couldn’t identify the language of the document in 4% of the documents.

Number of replicas	Number of Contents	% of contents
0	2462490	90.0
1	205882	7.5
2	33468	1.2
3	12814	0.5
4	6086	0.2
5	5272	0.2
6-10	6453	0.2
11-100	2318	0.1
101-1000	154	0.0
>1000	5	0.0
Total	2734942	100.0

Table VI. Distribution of contents with replicas.

6.5 Meta-tags

We studied the usage of two important meta-tags supported by HTML: *description* and *keywords* [W3C 1999]. The description meta-tag provides a description of the page's content and the meta-tag keywords provides a set of keywords that search engines may present as a result of a search. We found that just 17% of the pages had the meta-tag description and that, among these, the usage of this meta-tag doesn't seem to be correct. We found only 44000 distinct values for 555000 description meta-tags. This means that 92% of the texts of the descriptions were repeated elsewhere. We identified a set of causes for this situation:

- The meta tag value is a default text inserted by a publishing tool;
- The publisher repeated the same text in all the pages of its site, although they are different;
- There are replicated pages on the web.

The keywords meta-tag is present in 18% of the pages. A deeper analysis revealed that 91% of the pages that have the description meta-tag also had the keywords meta-tag. O'Neill et al. showed that the usage of meta-tags on the global web has been increasing in the past years. In 2002, 70% of the pages included meta-data [O'Neill et al. 2003]. Although our results focus only on two of the most popular meta-tags, we believe that the usage of meta-tags on the Portuguese web is much less frequent than on the global web.

The titles of the web pages aren't very descriptive either. There were over 600000 distinct titles for 3.1 million pages. The main reason we found for this observation, is that the title of the site's host page is used as the title for all the pages in the site in most cases.

7. WEB STRUCTURE

7.1 Content Replication

To detect contents replication, we compared the MD5 digest [Rivest 1992] obtained for each document. We identified 2734942 distinct contents. Table VI presents the replication distribution. We found that 15.5% of the contents were referenced

	# links	URL
1	3540	cpan.dei.uc.pt/modules/00modlist.long.html
2	2425	ftp.ist.utl.pt/pub/rfc/
3	2309	homepage.oninet.pt/095mad/bookmarks_on_mypage.html
4	1632	www.fis.uc.pt/bbsoft/bbhtm/mnusbib3.htm
5	1621	cpan.dei.uc.pt/authors/00whois.html
6	1532	www.fba.ul.pt/links4.html
7	1458	boa.oa.pt/bbsoft2/bbhtm/mnusbib3.htm
8	1346	www.esec-canecas.rcts.pt/Educacao/Escolas.htm
9	1282	pisco.cii.fc.ul.pt/nobre/hyt/bookmarks.html
10	1181	www.fpce.uc.pt/pessoais/rpaixao/9.htm

Table VII. The 10 documents with highest number of outgoing links.

by several distinct URLs (replicas). Mogul identified only 5% of replicas when analyzing a client trace from WebTV [Kelly and Mogul 2002]. We believe that the difference between our results, is due to the distinct methodologies adopted in the experiences. A crawl-based approach analyzes all the documents available on the web, while a client trace permits to analyze only the contents accessed by users. Most of the contents (90%) are unique and 7.5% had exactly one replica. Contents replicated more than 1000 times are very rare. However, they were the cause of 13146 downloads for just 5 distinct contents. These situations are pathological for web crawlers and also tend to bias the collection statistics. We observed these 5 cases and concluded that they were all caused by mal-functioning web servers, which always return the same error page for all the requests. Our measures against these traps (see Section 3.1) failed because all the links to documents with error messages were correctly extracted. When the crawler finally identified the trap, it already had numerous URLs to crawl, even though it had stopped inserting new links.

Our measurements indicate that 42% of the replicas are duplicates of contents hosted on the same site; 60% are duplicates of contents hosted on a different site; 2% are duplicates of contents hosted both in the same site and in another site.

7.2 Link Structure

The link structure of the web can be represented as a graph. Nodes represent URLs and edges represent hypertext links. We restricted our analysis to links between distinct sites, originated on web pages hosted under the .PT domain and targeted to one of the accepted TLDs. The links to URLs that evidenced that the referenced content was not of one of the accepted types, were excluded (e.g. URLs where the file part has a .jpg extension were not collected).

We observed that most of the web pages (95%), didn't link to another Portuguese site (Figure 10). On average a web page has 0.23 links to documents on another site. This is not a surprising result, since links are usually internal to the site. However, we also found pages rich in outgoing links (Table VII).

We found that 66% of the links didn't point to a document on the Portuguese web. This measure was made under the assumption that all the URLs hosted outside the .PT domain for which we couldn't determine the language were considered as being

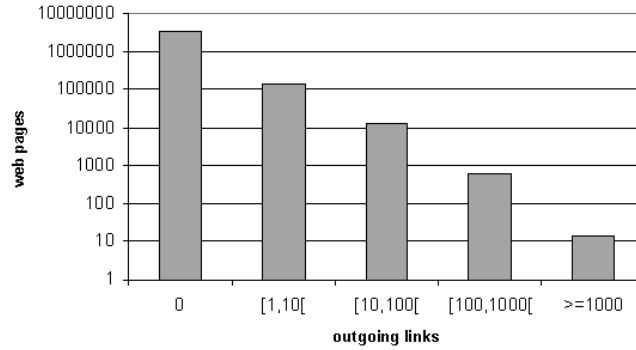


Fig. 10. Distribution of the number of outgoing links per web page.

	# incoming links	URL
1	6862	www.fccn.pt/
2	754	clanhosted.clix.pt
3	688	www.sapo.pt
4	606	www.publico.pt
5	522	www.infocid.pt
6	448	paginasbrancas.pt
7	423	www.dn.pt
8	413	www.sapo.pt/
9	361	security.vianetworks.pt
10	350	www.uminho.pt

Table VIII. The 10 documents with highest number of incoming links.

outside the Portuguese web. 40% of the links pointed to the host page of a site. We found that 3189710 documents (89%), were not referenced by a link originated in another Portuguese site. As observed on the global web, here most links tend to point to a small set of pages (Figure 11).

7.2.1 Document importance. The importance of a document can be determined through the analysis of the web graph. In order to achieve an meaningful ranking of the relative importance of documents, we handled links to replicas and HTTP redirects differently:

- Links to replicas cause the splitting of the number of links to a document among the several URLs that refer it. In the presence of replicated contents, we elected the document with the smallest URL as the common reference. We then erased the replicated pages from the web graph and re-targeted the links to the replicas to the URL used as common reference.
- HTTP redirects are almost invisible to web surfers. Involuntarily, publishers link to the URL of the redirect instead of the URL of the document. This causes a

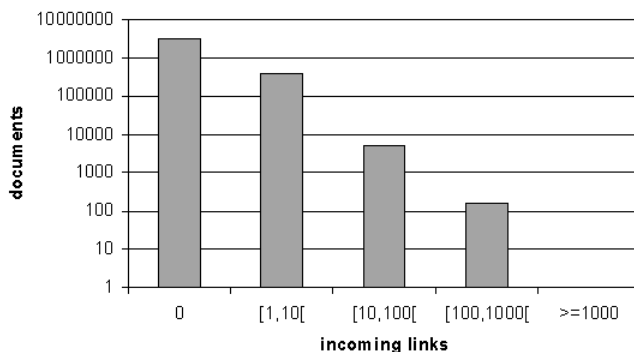


Fig. 11. Distribution of the number of incoming links per document.

	# incoming links	Site
1	7109	www.fccn.pt
2	1881	br.weather.com
3	1617	images.clix.pt
4	1601	www.sapo.pt
5	1481	www.clinicaviva.pt
6	794	www.depp.msst.gov.pt
7	777	www.infocid.pt
8	721	www.fct.mct.pt
9	652	ultimahora.publico.pt
10	615	www.miau.pt

Table IX. The 10 sites that received most incoming links

split in the number of links between the redirects and the document. We followed each redirect until we found a non-redirect URL. Then we replaced the redirect nodes in the graph by the correspondent non-redirect URLs.

Table VIII presents the 10 documents that received most incoming links on the web graph obtained after the above modifications to handle replicas and redirects were applied.

Despite our efforts to eliminate pathological situations in the analyzed graph, we still observe some anomalies on the most ranked document lists. In positions 3 and 8 of Table VIII, the number of incoming links was spread among 2 different URLs, although we know that they both refer to the same document. The problem was that we identified the two URLs by their string representation and between the crawl of the first and the second URL, the content referenced by them changed. Sometimes the change on the content is very small. In our second example, the change was just a link to an advertisement.

	# users	Site
1	779000	www.sapo.pt
2	580000	www.microsoft.com
3	560000	pesquisa.sapo.pt
4	548000	loginnet.passport.com
5	540000	www.clix.pt
6	538000	www.google.pt
7	480000	www.geocities.com
8	477000	login.passport.net
9	471000	www.terravista.pt
10	463000	www.iol.pt
11	408000	windowsupdate.microsoft.com
12	405000	v4.windowsupdate.microsoft.com
13	380000	www.msn.com
14	311000	pesquisa.clix.pt
15	290000	ww2.hpg.ig.com.br
16	247000	webmail.iol.pt
17	244000	www.mytmn.pt
18	227000	webmail.sapo.pt
19	224000	www.aeiou.pt
20	223000	www.google.com
21	219000	www.cidadebcp.pt
22	215000	planeta.clix.pt
23	203000	www.yahoo.com
24	202000	webmail.clix.pt
25	193000	caixadirecta.cgd.pt
26	191000	netcabo.sapo.pt
27	189000	dossieriraque.clix.pt
28	185000	tsf.sapo.pt
29	185000	www.cgd.pt
30	182000	login.passport.com
31	181000	www.msn.com.br
32	180000	bandalarga.netcabo.pt
33	177000	www.dgci.gov.pt
34	173000	www.abola.pt
35	171000	auth.clix.pt
36	165000	pwp.netcabo.pt
37	159000	www.tvi.iol.pt
38	156000	netbi.sapo.pt
39	156000	www.record.pt
40	156000	download.com.com

Table X. The 40 most accessed sites by Portuguese users. (courtesy 2003 Markttest Lda).

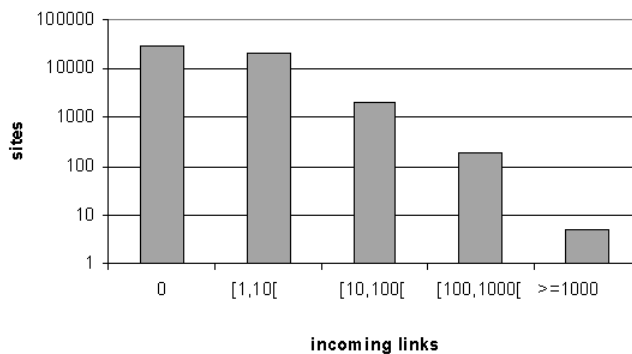


Fig. 12. Distribution of the number of incoming links per site.

7.2.2 Site importance. The importance of a site can be derived from the total number of incoming links. An highly ranked site might not host highly ranked documents. For instance, some online newspapers receive a large number of incoming links to many distinct news pages, but as news are interesting and are in many cases available for only a short period of time, they never get to be highly ranked documents.

Broder et al. analyzed the graph structure of the web through 2 large crawls of 200 million pages each [Broder et al. 2000]. They considered each page as a node and each hypertext link as an edge on the graph. They found that 91% of the pages were reachable from one another by following either forward or backward links after computing an algorithm that finds weak components in the graph. Our study followed a different methodology.

We considered only the links between distinct sites and didn't detect weak components in the graph. We generated a graph where each Portuguese site is a node and each link between documents on two different sites an edge.

We analyzed the graph as being undirected by following links in both directions and found that 73% of the sites connected to another site. This result contrasts with the one obtained by Broder (91%) and shows that the connectivity of the graph decreased on a smaller partition of the web, such as the Portuguese web.

Then we analyzed the graph as being directed, following links only in their real direction. We found that only 45% of the sites were reachable from one another site, which leaves us with a majority of sites (55%) that are never linked (orphan sites).

Figure 12 presents the distribution of the number of incoming links per site.

Table IX presents the 10 Portuguese sites that received most incoming links.

We obtained a list of 495 selected sites, accessed from the homes of a panel of Portuguese users, during the period we performed the crawl [Marktest 2003]. Table X presents the 40 sites that received most distinct users. We can observe that the majority (27) of these popular sites are hosted under the .PT domain as we assumed. We noticed a high number of accesses to sites that are automatically

accessed by tools. For instance, when a user types an URL of a site that is not found, Internet Explorer automatically redirects his request to `auto.search.msn.com` by default. These sites appear as overrated in usage statistics. We noticed that 50% of the sites accessed were part of the Portuguese web. We found a correlation of 0.527 between the number of users and links to the Portuguese sites. This shows that the most linked sites are also usually more visited by the users of this community web.

At first sight, it's surprising that the site `www.fccn.pt`, which occupies the first position in Table IX, is not present in the list of the 495 sites accessed by the users. A deeper analysis revealed that 96% of the links to the FCCN (National Foundation for Scientific Computing) site were originated on sites hosted under the `.RCTS.PT` domain and almost all of them (99%) pointed to the host page (`www.fccn.pt/`). The RCTS network (Network for Science, Technology and Society) is also managed by FCCN. It is composed by over 11000 sites from several public institutions, specially schools, hosted under the `.RCTS.PT` domain. We found that FCCN automatically generated a site on the RCTS network for every school in the country, initially composed by a single web page containing its address, e-mail and a link to `www.fccn.pt/`. The content of these sites was supposed to be replaced by contents produced by the schools, but in most cases this didn't happen. As a result, the default site prevailed, generating a high number of links to the FCCN site from other sites.

8. RELATED WORK

Web characterization has been done from different perspectives through the years almost since the beginning of the web [Pitkow 1998]. The Web Characterization Project has been a great contributor for research in web characterization [OCLC 2003; O'Neill et al. 2003].

Najork and Heydon performed a large scale web crawl from which they gathered statistics regarding the outcome of download attempts, distribution of types, size of the documents and replication [Najork and Heydon 2001]. They also witnessed that the distribution of pages over web servers follows a Zipfian distribution. Lawrence and Giles studied the accessibility of information on the web and drew conclusions about the size, extracted text and usage of meta-data in HTML pages [Lawrence and Giles 1999].

Boldi et al. studied the structural properties of the African web analyzing HTTP header fields and contents of HTML pages [Boldi et al. 2002] and Punpiti et al. presented quantitative measurements and analyses of documents hosted under the `.th` domain [Punpiti 2000].

Replication on the web has been studied in several works, through the syntactically clustering of documents [Broder et al. 1997], the study of the existence of near-replicas on the web [Shivakumar and Garcia-Molina 1998] and different levels of duplication between hosts and mechanisms to detect them [Bharat and Broder 1999]. The study of gateway and proxy traces also found replication on the web and identified that a few web servers are responsible for most of the duplicates [Douglis et al. 1997; Mogul 1999a]. A large client trace gathered from the WebTV networks evidenced the existence of URL aliasing and its implications to web caching systems [Kelly and Mogul 2002].

On language analysis, the authors propose a technique for estimating the size of language-specific corpus and used it to estimate the usage of English and non-English language on the WWW [Grefenstette and Nioche 2000]. Funredes presented a study on the presence of Latin languages on the web [Funredes 2001]. Aires et al. measured the web written in the Portuguese language [Aires and Santos 2002].

The notion of hostgraph and connectivity of web sites and country domains was presented in [Bharat et al. 2001].

A first effort to characterize the Portuguese web, defined a set of metrics to describe the web within the RCCN network (network that connected several Portuguese academic institutions) [Nicolau et al. 1997]. The Netcensus project aims at periodically collecting statistics regarding all type of files hosted under the .PT domain [Silva et al. 2002a; 2002b]. In our previous work, we presented a system for managing the deposit of digital publications and characterized a restricted set of Portuguese online publications, exposing the most common formats and file sizes [Noronha et al. 2001].

The statistics we gathered are sometimes significantly different from those presented in the bibliography. This is not a surprising result, since they are based in different and heterogeneous partitions of the web, using distinct methodologies and obtained in different dates.

9. CONCLUSION AND FUTURE WORK

This paper described our work in identifying, collecting and characterizing the Portuguese web. We propose defining criteria that cover this web with high precision, and is simultaneously easy to configure, when setting-up the harvesting policies on a crawler.

We observed that most of the sites are small virtual hosts under the .PT domain. The number of sites under construction is very high. The use of appropriate or descriptive meta-tags is still insignificant on the Portuguese web.

We identified situations on the web that may bias the results and proposed solutions, showing that web characterization depends on the used crawling technology.

This study is interesting to others who need to characterize community Webs and may help in the design of software systems that operate over the web. Web archivers can better estimate necessary resources and delimit partitions interesting for archival. Web proxies can be more accurately configured by administrators, crawlers can be improved through the definition of adequate architectures and crawling policies of web search engines can be used to improve their coverage of the web, leading to better search results.

As future work, we will extend the characterization of the Portuguese web to other MIME types and gather new metrics that would enable us to monitor the evolution of the web and its linkage structure. We also intend to improve the crawler performance so that statistics can be gathered in a shorter period of time. A major issue to be studied in the future is to define a more accurate and efficient definition of the Portuguese web. The current definition demands the download of large numbers of documents hosted outside the .PT domain, to identify a very small percentage written in Portuguese. This is highly inefficient, and makes it difficult for us to tell the documents of interest to our application domain from

others. The difficulty is particularly high for sites hosted under general purpose TLDs. In future work we intend to combine crawling policies with the usage of geographical tools such as Gtrace [Periakaruppan and Nemeth] in order to obtain a more precise definition.

10. ACKNOWLEDGEMENTS

We thank Miguel Costa and Bruno Martins for the discussions and development of software components that we used to extract the results presented in this paper. We thank Marktest for the access to the Netpanel statistics.

This study was partially supported by the FCCN-Fundação para a Computação Científica Nacional, FCT-Fundação para a Ciência e Tecnologia, under grants POSI/ SRI/ 40193/ 2001 (XMLBase project) and SFRH/ BD/ 11062/ 2002 (scholarship).

REFERENCES

- AIRES, R. AND SANTOS, D. 2002. Measuring the web in Portuguese. In *Euroweb 2002 Conference*, B. Matthews, B. Hopgood, and M. Wilson, Eds. Oxford, UK, 198–199.
- ALBERTSEN, K. 2003. The paradigm web harvesting environment. In *Proceedings of 3rd ECDL Workshop on Web Archives*. Trondheim, Norway.
- BARR, D. 1996. *Common DNS Operational and Configuration Errors*.
- BHARAT, K. AND BRODER, A. 1999. Mirror, mirror on the web: a study of host pairs with replicated content. In *Proceedings of the Eighth International Conference on World Wide Web*. Elsevier North-Holland, Inc., 1579–1590.
- BHARAT, K., CHANG, B.-W., HENZINGER, M. R., AND RUHL, M. 2001. Who links to whom: Mining linkage between web sites. In *Proceedings of the 2001 IEEE International Conference on Data Mining*. IEEE Computer Society, 51–58.
- BOLDI, P., CODENOTTI, B., SANTINI, M., AND VIGNA, S. 2002. Structural properties of the African web. In *Proceedings of the 11th International World Wide Web Conference*. Honolulu, Hawaii.
- BRIN, S. AND PAGE, L. 1998. The anatomy of a large-scale hypertextual web search engine. *Computer Networks and ISDN Systems* 30, 1–7, 107–117.
- BRODER, A., KUMAR, R., MAGHOUL, F., RAGHAVAN, P., RAJAGOPALAN, S., STATA, R., TOMKINS, A., AND WIENER, J. 2000. Graph structure in the web. In *Proceedings of the 9th International World Wide Web Conference on Computer networks*. North-Holland Publishing Co., 309–320.
- BRODER, A. Z., GLASSMAN, S. C., MANASSE, M. S., AND ZWEIG, G. 1997. Syntactic clustering of the web. In *Proceedings of the Sixth International conference on World Wide Web*. Elsevier Science Publishers Ltd., 1157–1166.
- CAVNAR, W. AND TRENKLE, J. 1994. N-gram-based text categorization. In *Third Annual Symposium on Document Analysis and Information Retrieval*. 161–175.
- CENTER, H. S. D. 2003. Geo targeting ip address to country city region isp latitude longitude database for internet developers - ip2location. <http://www.ip2location.com/>.
- CHO, J. AND GARCIA-MOLINA, H. 2000. The evolution of the web and implications for an incremental crawler. In *VLDB 2000, Proceedings of 26th International Conference on Very Large Data Bases, September 10-14*. 200–209.
- CHO, J., GARCIA-MOLINA, H., AND PAGE, L. 1998. Efficient crawling through URL ordering. *Computer Networks and ISDN Systems* 30, 1–7, 161–172.
- DAVIS, C., VIXIE, P., GOODWIN, T., AND DICKINSON, I. 1996. *A Means for Expressing Location Information in the Domain Name System*.
- DAY, M. 2003. Collecting and preserving the world wide web. http://www.jisc.ac.uk/uploaded_documents/archiving_feasibility.pdf.
- DOUGLIS, F., FELDMANN, A., KRISHNAMURTHY, B., AND MOGUL, J. C. 1997. Rate of change and other metrics: a live study of the world wide web. In *USENIX Symposium on Internet Technologies and Systems*.

- FETTERLY, D., MANASSE, M., NAJORK, M., AND WIENER, J. L. 2003. A large-scale study of the evolution of web pages. In *Proceedings of the 12th International World Wide Web Conference*. Budapest.
- FLAKE, G., LAWRENCE, S., AND GILES, C. L. 2000. Efficient identification of web communities. In *Sixth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. Boston, MA, 150–160.
- FUNREDES. 2001. The place of latin languages on the internet. http://www.funredes.org/LC/english/L5/L5index_english.html.
- GIBSON, D., KLEINBERG, J. M., AND RAGHAVAN, P. 1998. Inferring web communities from link topology. In *Proceedings of the 9th ACM Conference on Hypertext and Hypermedia*. Pittsburgh, Pennsylvania, 225–234.
- GOMES, D. 2003. Viúva negra. www.tumba.pt/english/crawler.html.
- GOOGLE. 2003. Google web search features. www.google.com/help/features.html#link.
- GREFENSTETTE, G. AND NIOCHE, J. 2000. Estimation of english and non-english language use on the WWW. In *Proceedings of RIAO'2000, Content-Based Multimedia Information Access*. Paris, 237–246.
- HARRENTIEN, K., STAHL, M. K., AND FEINLER, E. J. 1985. *NICNAME/WHOIS*.
- HENZINGER, M. 2003. Algorithmic challenges in web search engines. *Journal of Internet Mathematics* 1, 1, 115–126.
- HEYDON, A. AND NAJORK, M. 1999. Mercator: A scalable, extensible web crawler. *World Wide Web* 2, 4, 219–229.
- KELLY, T. AND MOGUL, J. 2002. Aliasing on the world wide web: Prevalence and performance implications. In *Proceedings of the 11th International World Wide Web Conference*. Honolulu, Hawaii.
- LAWRENCE, S. AND GILES, C. L. 1999. Accessibility of information on the web. *Nature* 400, 107–109.
- LEUNG, S.-T. A., PERL, S. E., STATA, R., AND WIENER, J. L. 2001. Towards web-scale web archeology. Research Report 174, Compaq Research Center, Paolo Alto CA. September.
- LLC, M. 2003. Maxmind: How to locate your internet visitors geotargeting ip address to country state city isp organization latitude longitude. <http://www.maxmind.com/>.
- MARKTEST. 2003. Netpanel. netpanel.marktest.pt/.
- MOGUL, J. 1999a. A trace-based analysis of duplicate suppression in HTTP. Technical Report 99/2, Compaq Computer Corporation Western Research Laboratory. November.
- MOGUL, J. 1999b. Errors in timestamp-based HTTP header values. Research Report 99/3, Compaq Computer Corporation Western Research Laboratory. December.
- NAJORK, M. AND HEYDON, A. 2001. On high-performance web crawling. Src research report, Compaq Systems Research Center.
- NETCRAFT LTD. 2004. Netcraft: April 2003 archives. <http://news.netcraft.com/archives/2003/04/index.html>.
- NICOLAU, M. J., MACEDO, J., AND COSTA, A. 1997. Caracterização da informação WWW na RCCN. Tech. rep., Universidade do Minho.
- NORONHA, N., CAMPOS, J. P., GOMES, D., SILVA, M. J., AND BORBINHA, J. 2001. A deposit for digital collections. In *Proc. 5th European Conf. Research and Advanced Technology for Digital Libraries, ECDL*. Springer-Verlag, 200–212.
- OCLC. 2003. Web characterization. <http://wcp.oclc.org/>.
- O'NEILL, E. T. 1999. Web sites: Concepts, issues, and definitions. <http://wcp.oclc.org/pubs/rn1-websites.html>.
- O'NEILL, E. T., LAVOIE, B. F., AND BENNETT, R. 2003. Trends in the evolution of the public web. *D-Lib Magazine* 9, 4 (April).
- OVERTURE SERVICES, I. 2003. Alltheweb.com: Frequently asked questions - url investigator. www.alltheweb.com/help/faqs/url_investigator.
- PERIAKARUPPAN, R. AND NEMETH, E. GTrace: A graphical traceroute tool. 69–78.

- PITKOW, J. E. 1998. Summary of WWW characterizations. *Computer Networks and ISDN Systems* 30, 1–7, 551–558.
- POSTEL, J. 1994. *Domain Name System Structure and Delegation*.
- PUNPITI, S. S. 2000. Measuring and analysis of the Thai world wide web. In *Proceedings of the Asia Pacific Advance Network*. 225–230.
- RIVEST, R. 1992. *RFC 1321 - The MD5 Message-Digest Algorithm*.
- SHIVAKUMAR AND GARCIA-MOLINA. 1998. Finding near-replicas of documents on the web. In *Workshop on Web Databases (WebDB'98)*. LNCS.
- SILVA, L. O., MACEDO, J., COSTA, A., BELO, O., AND SANTOS, A. 2002a. Netcensus: Medição da evolução dos conteúdos na web. Tech. rep., Departamento de Informática, Universidade do Minho.
- SILVA, L. O., MACEDO, J., COSTA, A., BELO, O., AND SANTOS, A. 2002b. Obtenção de estatísticas do www em Portugal. Tech. rep., OCT and DI, Universidade do Minho.
- SILVA, M. J. 2003. The case for a portuguese web search engine. In *Proceedings of IADIS International Conference WWW/Internet 2003*. Algarve, Portugal.
- SPINELLIS, D. 2003. The decay and failures of web references. *Communications of the ACM* 46, 1, 71–77.
- W3C. 1999. HTML 4.01 specification. <http://www.w3.org/TR/html401/>.
- W3C. 1999. Web characterization terminology and definitions sheet. <http://www.w3.org/1999/05/WCA-terms/>.
- WEBB, C. 2000. Towards a preserved national collection of selected Australian digital publications. In *Proceedings of Preservation 2000 Conference*. York, UK.
- WILLS, C. E. AND MIKHAILOV, M. 1999. Towards a better understanding of Web resources and server responses for improved caching. *Computer Networks (Amsterdam, Netherlands: 1999)* 31, 11–16, 1231–1243.
- ZABICKA, P. 2003. Archiving the Czech web: issues and challenges. In *Proceedings of 3rd ECDL Workshop on Web Archives*. Trondheim, Norway.
- ZOOK, M. 2000. Internet metrics: Using host and domain counts to map the internet.