

This paper was unfortunately omitted from the print version of the June 2008 issue of Sigir Forum.  
Please cite as SIGIR Forum June 2008, Volume 42 Number 1, pp 78-83

# A Large Time-Aware Web Graph

Paolo Boldi                      Massimo Santini  
*boldi@dsi.unimi.it*            *santini@dsi.unimi.it*

Sebastiano Vigna  
*vigna@dsi.unimi.it*

Dipartimento di Scienze dell'Informazione  
Università di Milano, Italy

## Abstract

We describe the techniques developed to gather and distribute in a highly compressed, yet accessible, form a series of twelve snapshot of the .uk web domain. *Ad hoc* compression techniques made it possible to store the twelve snapshots using just 1.9 bits per link, with constant-time access to temporal information. Our collection makes it possible to study the temporal evolution link-based scores (e.g., PageRank), the growth of online communities, and in general time-dependent phenomena related to the link structure.

## 1 Introduction

By now, several sources provide accessible snapshots of web data. The Stanford WebBase project, for instance, provides hundreds of Terabytes of such data. In the last years, the LAW (Laboratory for Web Algorithmics) focused its efforts on web graphs instead.

Within the DELIS project, interest rose about *temporal* link analysis, that is, studying how the web-graph nodes and arcs evolves in time. Since without proper bounds the amount of information required by this activity is staggering, we decided to concentrate our efforts on gathering twelve monthly 100 Mpages snapshots of the .uk domain and store them in a format that would make temporal link information accessible on a standard workstation.

The basis of our work is WebGraph [2], a framework for web graph compression that currently provides the best compression available (in terms of bits per link) and whose data can be accessed using free Java or C++ code [6]. One of the main challenges of this work was to extend WebGraph so that it would compress efficiently labels representing temporal information. Moreover, we wanted to extend the standard WebGraph flexible approach, that

---

makes it always possible to load data into main memory or access it in offline form (with a performance drop, of course).

## 2 Gathering the Snapshots

The snapshots have been taken at the start of each month, during a period of 7 – 10 days, using the bandwidth provided by the Università degli Studi di Milano and a cluster of PCs that has been funded by the DELIS project. Some basic information about the snapshots is shown in Table 1.

**Crawling parameters.** As in any limited-size crawl, it is essential to define the crawl parameters (the stopping criterion is clearly that of reaching about 100 Mpages, without counting duplicates). We highlight the main features:

**Crawl policy.** We use UbiCrawler’s [1] built-in per-host breadth-first visit. A number of threads scan in parallel distinct hosts, and newly discovered URLs are added to a queue. When a thread completes its visit, it extracts from the queue the first URL whose host has an IP address that is not currently being visited, and starts visiting that host in breadth-first fashion.

**Seed.** The seed is a large (190 000 elements) set of URLs obtained from the Open Directory Project. The reason for such a large seed is that of making the crawl more stable and repeatable, and reduce the amount of spam (as links in the Open Directory Project are judged by humans).

**Maximum number of pages per host.** We limited each host to a maximum of 50 000 pages. This guarantees that we shall crawl at least 2 000 hosts, and limits the impact of web traps and database-driven sites.

**Maximum inter-host depth.** We do not delve more than 16 levels in a host. The main reason for a limit in depth is avoiding traps and also badly configured 404 pages, which sometimes generate an infinite number of links by prefix buildup.

**URL normalisation.** URLs are normalised following the strategy explained in the BURL<sup>1</sup> Java class. We apply all safe normalisations, escape all illegal characters, and treat in a special way square brackets as they are ubiquitously (although erroneously) used in an unescaped form.

**Duplicate detection.** Many pages are duplicates, and to detect their presence we maintain a set of 64-bit fingerprints obtained after stripping attributes (of HTML elements) and other non-relevant parts of the page. When a duplicate is detected we just store a pointer to the original page. About 25% of the overall pages happen to be duplicates.

---

<sup>1</sup>The class is available in bundle with the LAW software, downloadable at <http://law.dsi.unimi.it/>.

---

	Pages	Size (GB)	GZip'd Size (GB)
June	112 386 763	1 893	402
July	136 956 559	2 287	477
August	141 395 895	2 424	507
September	148 965 298	2 756	546
October	129 558 491	2 336	478
November	150 146 132	2 637	546
December	144 489 446	2 552	525
January	151 578 113	2 651	553
February	153 966 540	2 692	564
March	151 427 461	2 568	545
April	150 606 689	2 700	559
May	150 054 551	2 658	556

	Nodes	Arcs	Size (GB)	bit/arc
June	80 644 902	2 481 281 617	0.89	3.07
July	96 395 298	3 030 665 444	1.16	3.30
August	100 751 978	3 250 153 746	1.23	3.25
September	106 288 541	3 871 625 613	1.32	2.93
October	93 463 772	3 130 910 405	1.03	2.83
November	106 783 458	3 479 400 938	1.16	2.86
December	103 098 631	3 768 836 665	1.34	2.77
January	108 563 230	3 929 837 236	1.38	2.72
February	110 123 614	3 944 932 566	1.39	2.74
March	107 565 084	3 642 701 825	1.34	2.84
April	106 867 191	3 790 305 474	1.36	2.79
May	105 896 555	3 738 733 648	1.30	2.69

Table 1: Per-snapshot full-text and web-graph stats.





