

Self-similarity in the web

Stephen Dill Ravi Kumar Kevin McCurley
Sridhar Rajagopalan D. Sivakumar Andrew Tomkins

February 26, 2001

Submission to the 27th International Conference on Very Large Databases (VLDB 2001)

Category: Research

Topic Area: Databases and database services in new context - Internet and the WWW

Contact Author:

Ravi Kumar

IBM Almaden Research Center

650 Harry Road

ravi@almaden.ibm.com

408-927-1885

Self-similarity in the web*

Stephen Dill Ravi Kumar Kevin McCurley Sridhar Rajagopalan
D. Sivakumar Andrew Tomkins

February 26, 2001

Abstract

Algorithmic tools for searching and mining the web are becoming increasingly sophisticated and vital. In this context, algorithms which use and exploit structural information about the web perform better than generic methods in both efficiency and reliability.

We present an extensive characterization of the graph structure of the web, with a view to enabling high-performance applications that make use of this structure. In particular, we show that the web emerges as the outcome of a number of essentially independent stochastic processes that evolve at various scales. A striking consequence of this scale invariance is that the structure of the web is “fractal” — cohesive sub-regions display the same characteristics as the web at large. An understanding of this underlying fractal nature is therefore applicable to designing data services across multiple domains and scales.

We describe potential applications of this line of research to optimized algorithm design for web-scale data analysis.

*IBM Almaden Research Center, 650 Harry Road, San Jose, CA 95120.

1 Introduction

As the size of the web grows exponentially, data services on the web are becoming increasingly complex and challenging tasks. These include both basic services such as searching and finding related pages, and advanced applications such as web-scale data mining, community extraction, constructions of indices, taxonomies, and vertical portals. Applications are beginning to emerge that are required to operate at various points on the “petabyte curve” — billions of web pages that each have megabytes of data, tens of millions of users in a peer-to-peer setting each with several gigabytes of data, etc. The upshot of the rate and diversity of this growth is that data service applications for collections of hyperlinked documents need to be efficient and effective at several scales of operation. As we will show, a form of “scale invariance” exists on the web that allows simplification of this multi-scale data service design problem.

The first natural approach to the wide range of analysis problems emerging in this new domain is to develop a general query language to the web. There have been a number of proposals along these lines [34, 6, 43]. Further, various advanced mining operations have been developed in this model using a web-specific query language like those described above, or a traditional database encapsulating some domain knowledge into table layout and careful construction of SQL programs [18, 42, 8].

However, these applications are particularly successful precisely when they take advantage of the special structure of the document collections and the hyperlink references among them. An early example of this phenomenon in the marketplace is the paradigm shift witnessed in search applications — ranking schemes for web pages that were based on link analysis [26, 12] proved to be vastly superior to the more traditional text-based ones.

The success of these specialized approaches naturally led researchers to seek a finer understanding of the hyperlinked structure of the web. Broadly, there are two (very related) lines of research that have emerged. The first one is more theoretical and is concerned with proposing stochastic models that explain the hyperlink structure of the web [27, 7, 1]. The second line of research [13, 7, 3, 28] is more empirical; new experiments are conducted that either validate or refine existing models.

There are several driving applications that motivate (and are motivated by) a better understanding of the neighborhood structure on the web. In particular, the “second generation” of data service applications on the web — including advanced search applications [16, 17, 10], browsing and information foraging [14, 39, 15, 40, 19], community extraction [28], taxonomy construction [30, 29] — have all taken tremendous advantage of knowledge about the hyperlink structure of the web. As just one example, let us mention the community extraction algorithm of [28]. In this algorithm, a characterization of degree sequences within web-page neighborhoods allowed the development and analysis of efficient pruning algorithms for a sub-graph enumeration problem that is in general intractable.

Even more recently, new algorithms have been developed to benefit from structural information about the web. Arasu *et al.* [5] have shown how to take advantage of the macroscopic “bow-tie” structure of the web [13] to design an efficient algorithmic partitioning method for certain eigenvector computations; these are the key to the successful search algorithms of [26, 12], and to popular database indexing methods such as latent semantic indexing [20, 36]. Adler and Mitzenmacher [4] have shown how the random graph characterizations of the web given in [27] can be used to construct very effective strategies to compress the web graph.

1.1 Our results

In this paper, we present a much more refined characterization of the structure of the web. Specifically, we present evidence that the web emerges as the outcome of a number of essentially independent stochastic processes that evolve at various scales, all roughly following the model of [27]. A striking consequence of this is that the web is a “fractal” — each thematically unified region displays the same characteristics as the web at large. This implies the following useful corollary:

To design efficient algorithms for data services at various scales on the web (vertical portals pertaining to a theme, corporate intranets, etc.), it is sufficient (and perhaps necessary) to understand the structure that emerges from one fairly simple stochastic process.

We believe that this is a significant step in web algorithmics. For example, it shows that the sophisticated algorithms of [5, 4] are only the beginning, and the prospects are, in fact, much wider. We fully expect future data applications on the web to leverage this understanding.

Our characterization is based on two findings we report in this paper. Our first finding is an experimental result. We show that self-similarity holds for many different parameters, and also for many different approaches to defining varying scales of analysis. Our second finding is an interpretation of the experimental data. We show that, at various different scales, cohesive collections of web pages (for instances, pages on a site, or pages about a topic) mirror the structure of the web at large. Furthermore, if the web is decomposed into these cohesive collections, for a wide range of definitions of “cohesive,” the resulting collections are tightly and robustly connected via a *navigational backbone* that affords strong connectivity between the collections. This backbone ties together the collections of pages, but also ties together the many different and overlapping *decompositions* into cohesive collections, suggesting that committing to a single taxonomic breakdown of the web is neither necessary nor desirable. We now describe these two findings in more detail.

First, self-similarity in the web is pervasive and robust — it applies to a number of essentially independent measurements and regardless of the particular method used to extract a slice of the web. Second, we present a graph-theoretic interpretation of the first set of observations, which leads to a natural hierarchical characterization of the structure of the web interpreted as a graph. In our characterization, collections of web pages that share a common attribute (for instance, all the pages on a site, or all the pages about a particular topic) are structurally similar to the whole web. Furthermore, there is a *navigational backbone* to the web that provides tight and robust connections between these focused collections of pages.

1. Experimental findings. Our first finding, that self-similarity in the web is pervasive and appears in many unrelated contexts, is an experimental result. We explore a number of graph-theoretic and syntactic parameters. The set of parameters we consider is the following: indegree and outdegree distributions; strongly- and weakly- connected component sizes; bow-tie structure and community structure on the web graph; and population statistics for trees representing the URL namespace. We define these parameters formally below. We also consider a number of methods for decomposing the web into interesting subgraphs. The set of subgraphs we consider is the following: a large internet crawl; various subgraphs consisting of about 10% of the sites in the original crawl; 100 websites from the crawl each containing at least 10,000 pages; ten graphs, each consisting of every page containing a set of keywords (in which the ten keyword sets represent five broad topics and five sub-topics of the broad topics); a set of pages containing geographical references (e.g., phone numbers, zip codes, city names, etc.) to locations in the western United States; a graph representing the connectivity of web sites (rather than web pages); and a crawl of the IBM intranet.

We then consider each of the parameters described above, first for the entire collection, and then for each decomposition of the web into sub-collections. Self-similarity is manifest in the resulting measurements in two flavors. First, when we fix a collection or sub-collection and focus on the distribution of any parameter (such as the number of hyperlinks, number of connected components, etc.), we observe a Zipfian self-similarity within the pageset.¹ Namely, for any parameter x with distribution X , there is a constant c such that for all $t > 0$ and $a \geq 1$, $X(at) = a^c X(t)$.² Second, the phenomena (whether distributional or structural) that are manifest within a sub-collection are also observed (with essentially the same constants) in the entire collection, and more generally, in all sub-collections at all scales — from local websites to the web as a whole.

2. Interpretations. Our second finding is an interpretation of the experimental data. As mentioned above, the sub-collections we study are created to be cohesive clusters, rather than simply random sets of web pages. We will refer to them as *thematically unified clusters*, or simply TUCs. Each TUC has structure similar to the web as a whole. In particular, it has a Zipfian distribution over the parameters we study, strong navigability properties, and significant community and bow-tie structure (in a sense to be made explicit below).

Furthermore, we observe unexpectedly that the central regions of different TUCs are tightly and robustly connected together. These tight and robust inter-cluster linking patterns provide a *navigational backbone* for the web. By analogy, consider the problem of navigating from one physical address to another. A user might take a cab to the airport, take a flight to the appropriate destination city, and take a cab to the destination address. Analogously, navigation between TUCs is accomplished by traveling to the central core of a TUC, following the navigational backbone to the central core of the destination TUC, and finally navigating within the destination TUC to the correct page. We show that the self-similarity of the web graph, and its local and global structure, are alternate and equivalent ways of viewing this phenomenon.

1.2 Related prior work

Zipf-Pareto-Yule and Power laws.

Distributions with an inverse polynomial tail have been observed in a number of contexts. The earliest observations are due to Pareto [38] in the context of economic models. Subsequently, these statistical behaviors have been observed in the context of literary vocabulary [45], sociological models [46], and even oligonucleotide sequences [33], among others. Our focus is on the closely related power law distributions, defined on the positive integers, with the probability of the value i being proportional to i^{-k} for a small positive number k . Perhaps the first rigorous effort to define and analyze a model for power law distributions is due to Herbert Simon [41].

Recent work [30, 7] suggests that both the in- and the outdegrees of nodes on the web graph have power laws. The difference in scope in these two experiments is noteworthy. The first [30] examines a web crawl from 1997 due to Alexa, Inc., with a total of over 40 million nodes. The second [7] examines web pages from the University of Notre Dame domain `*.nd.edu` as well as a portion of the web reachable from 3 other URLs. This collection of findings already leads us to suspect the fractal-like structure of the web.

¹For more about the connection between Zipfian distributions and self-similarity, see Section 2.2 and [31].

²For example, the fraction of web pages that have k hyper-inlinks is proportional to $k^{-2.1}$.

Graph-theoretic methods.

Much recent work has addressed the web as a graph and applied algorithmic methods from graph theory in addressing a slew of search, retrieval, and mining problems on the web. The efficacy of these methods was already evident even in early local expansion techniques [14]. Since then, increasingly sophisticated techniques have been used; the incorporation of graph-theoretical methods with both classical and new methods that examine both context and content, and richer browsing paradigms have enhanced and validated the study and use of such methods. Following Botafogo and Shneiderman [14], the view that connected and strongly-connected components represent meaningful entities has become widely accepted.

Power laws and browsing behavior.

The power law phenomenon is not restricted to the web graph. For instance, [21] report very similar observations about the physical topology of the internet. Moreover, the power law characterizes not only the structure and organization of information and resources on the web, but also the way people use the web. Two lines of work are of particular interest to us here. (1) Web page access statistics, which can be easily obtained from server logs (but for caching effects) [22, 25, 2]. (2) User behavior, as measured by the number of times users at a single site access particular pages also enjoy power laws, as verified by instrumenting and inspecting logs from web caches, proxies, and clients [9, 32].

There is no direct evidence that browsing behavior and linkage statistics on the web graph are related in any fundamental way. However, making the assumption that linkage statistics directly determine the statistics of browsing has several interesting consequences. The Google search algorithm, for instance, is an example of this. Indeed, the view of PageRank put forth in [12] is that it puts a probability value on how easy (or difficult) it is to find particular pages by a browsing-like activity. Moreover, it is generally true (for instance, in the case of random graphs) that this probability value is closely related to the indegree of the page. In addition there is recent theoretical evidence [27, 41] suggesting that this relationship is deeper. In particular, if one assumes that the ease of finding a page is proportional to its graph-theoretic indegree, and that otherwise the process of evolution of the web as a graph is a random one, then power law distributions are a direct consequence. The resulting models, known as *copying* models for generating random graphs seem to correctly predict several other properties of the web graph as well.

2 Preliminaries

In this section we formalize our view of the web as a graph; here we ignore the text and other content in pages, focusing instead on the links between pages. In the terminology of graph theory [23], we refer to pages as *nodes*, and to links as *arcs*. In this framework, the web is a large graph containing over a billion nodes, and a few billion arcs.

2.1 Graphs and terminology

A *directed graph* consists of a set of *nodes*, denoted V and a set of *arcs*, denoted E . Each arc is an ordered pair of nodes (u, v) representing a directed connection from u to v . The *outdegree* of a node u is the number of distinct arcs $(u, v_1), \dots, (u, v_k)$ (i.e., the number of links from u), and the *indegree* is the number of distinct arcs $(v_1, u), \dots, (v_k, u)$ (i.e., the number of links to u). A path from node u to node v is a sequence of arcs $(u, u_1), (u_1, u_2), \dots, (u_k, v)$. One can follow such a sequence of arcs to “walk” through the graph from u to v . Note that a path from u to v does not imply a path from v to u . The *distance* from u to v

is one more than the smallest k for which such a path exists. If no path exists, the distance from u to v is defined to be infinity. If (u, v) is an arc, then the distance from u to v is 1. Given a graph (V, E) and a subset V' of the node set v , the *node-induced subgraph* (V', E') of (V, E) is defined by taking E' to be $\{(u, v) \in E \mid u, v \in V'\}$, i.e., the node-induced subgraph corresponding to some subset V' of the nodes contains only arcs that lie entirely within V' .

Given a directed graph, a *strongly connected component* of this graph is a set of nodes such that for any pair of nodes u and v in the set there is a path from u to v . In general, a directed graph may have one or many strong components. Any graph can be partitioned into a disjoint union of strong components. Given two strongly connected components, C_1 and C_2 , either there is a path from C_1 to C_2 or a path from C_2 to C_1 or neither, but not both. Let us denote the largest strongly component by *SCC*. Then, all other components can be classified with respect to the SCC in terms of whether they can reach, be reached from, or are independent of, the SCC. Following [13], we denote these components *IN*, *OUT*, and *OTHER* respectively. The SCC, flanked by the IN and OUT, figuratively forms a “bow-tie.”

A *weakly connected component* (*WCC*) of a graph is a set of nodes such that for any pair of nodes u and v in the set, there is a path from u to v if we disregard the directions of the arcs. Similar to strongly connected components, the graph can be partitioned into a disjoint union of weakly connected components. We denote the largest weakly connected component by *WCC*.

2.2 Zipf distributions and power laws

The power law distribution with parameter $a > 1$ is a distribution over the positive integers. Let X be a power law distributed random variable with parameter a . Then, the probability that $X = i$ is proportional to i^{-a} . The Zipf distribution is an interesting variant on the power law. The Zipf distribution is defined over any categorical-valued attribute (for instance, words of the English language). In the Zipf distribution, the probability of the i -th most likely attribute value is proportional to i^{-a} . Thus, the main distinction between these is in the nature of the domain from which the r.v. takes its values. A classic general technique for computing the parameter a characterizing the power law is due to Hill [24]. We will use Hill’s estimator as the quantitative measure of self-similarity.

While a variety of socio-economic phenomena have been observed to obey Zipf’s law, there is only a handful of stochastic models for these phenomena of which satisfying Zipf’s law is a consequence. Simon [41] was perhaps the first to propose a class of stochastic processes whose distribution functions follow the Zipf law [31]. Recently, new models have been proposed for modeling the evolution of the web graph [27]. These models predict that several interesting parameters of the web graph obey the Zipf law.

3 Experimental setup

3.1 Random subsets and TUCs

Since the average degree of the web graph is small, one should expect subgraphs induced by (even fairly large) random subsets of the nodes to be almost empty. Consider for instance a random sample of 1 million web pages (say out of a possible 1 billion pages). Consider now an arbitrary arc, say (a, b) . The probability that both endpoints of the arc are chosen in the random sample is about 1 in a million ($1/1000 * 1/1000$). Thus, the total expected number of arcs in the induced subgraph of these million nodes is about 8000, assuming an average degree of 8 for the web as a whole. Thus, it would be unreasonable to expect random

subgraphs of the web to contain any graph-theoretic structure. However, if the subgraphs chosen are not random, the situation could be (and is) different. In order to highlight this dichotomy, we introduce the notion of a *thematically unified cluster (TUC)*. A TUC is a cluster of webpages that share a common trait. In all instances we consider, these thematically unified clusters share a fairly syntactic trait. However, we do not wish to restrict our definition only to such instances. For instance, one could consider linkage-based concepts [43], and [39] as well. We now detail several instances of TUCs.

(1) *By content*: The premise that web content on any particular topic is also “local” in a graph-theoretic context has motivated some interesting earlier work [26, 30]. Thus, one should expect web pages that share subject matter to be more densely linked than random subsets of the web. If so, these graphs should display interesting morphological structure. Moreover, it is reasonable to expect this structure to represent interesting ways of further segmenting the topic.

The most naive method for judging content correlation is to simply look at a collection of webpages which share a small set of common keywords. To this end, we have generated 10 slices of the web, denoted henceforth as KEYWORD1, . . . , KEYWORD10. To determine whether a page belongs to a keyword set, we simply look for the keyword in the body of the document after simple pre-processing (removing tags, javascript, transform to lower case, etc.). The particular keyword sets we consider are shown in Tables 3 and 4 below. The terms in the first table correspond to mesoscopic subsets and the corresponding terms in the second table are microscopic subsets of the earlier ones.

(2) *By location*: Websites and intranets are logically consistent ways of partitioning the web. Thus, they are obvious candidates for TUCs. We look at intranets and particular websites to see what structures are represented at this level. We are interested in what features, if any, distinguish these two cases from each other and indeed from the web at large. Our observations here would help determine what special processing, if any, would be relevant in the context of an intranet. To this end, we have created TUCs consisting of the IBM intranet, denoted INTRANET henceforth, and 100 websites denoted SUBDOMAIN1, . . . , SUBDOMAIN100, each containing at least 10K pages.

(3) *By geographic location*: Geography is becoming increasingly evident in the web, with the growth in the number of local and small businesses represented on the web (restaurants, shows, housing information, and other local services) as well as local information websites such as `sidewalk.com`. We expect the recurrence of similar information structures at this level. We hope to understand more detail about overlaying geospatial information on top of the web. We have created a subset of the web based on geographic cues, denoted GEO henceforth. The subset contains pages that have geographical references (addresses, telephone numbers, and ZIP codes) to locations in the western United States. This was constructed through the use of databases for latitude–longitude information for telephone number area codes, prefixes, and postal zipcodes. Any page that contained a zipcode or telephone number was included if the reference was within a region bounded by Denver (Colorado) on the east and Nilolski (Alaska) on the west, Vancouver (British Columbia) on the north, and Brownsville (Texas) on the south.

To complete our study, we also define some additional graphs derived from the web. Strictly speaking, these are not TUCs. However, they can be derived from the web in a fairly straightforward manner. As it turns out, some of our most interesting observations about the web relates to the interplay between structure at the level of the TUCs and structure at the following levels. We define them now:

(4) *Random collections of websites*: We look at all the nodes that belong in a random collection of websites. We do this in order to understand the fine grained structure of the SCC, which is the navigational backbone of the web. Unlike random subgraphs of the web, random collections of websites exhibit interesting behaviors. First, the local arcs within a website ensure that there is fairly tight connectivity within each website. This allows the small number of additional intersite arcs to be far more useful than would be the case in a random subgraph. We have generated 7 such disjoint subsets. We denote these STREAM1, . . . , STREAM7.

(5) *Hostgraph*: The hostgraph contains a single node corresponding to each website (for instance `www.ibm.com` is represented by a single node), and has an arc between two nodes, whenever there is a page in the first website that points to a page in the second. The hostgraph is not a subgraph of the web graph, but it can be derived from it in a fairly straightforward manner, and more importantly, is relevant to understanding the structure of linkage at levels higher than that of a web page. In the following discussion, this graph is denoted by HOSTGRAPH.

3.2 Parameters studied

We study the following parameters:

(1) *Indegree distributions*: Recall that the indegree of a node is the number of arcs whose destination is that node. We consider the distribution of indegree over all nodes in a particular graph, and consider properties of that distribution. A sequence of papers [7, 3, 28, 13] has provided convincing evidence that indegree distributions follow the power law, and that the parameter a (called *indegree exponent*) is reliably around 2.1 (with little variation). We study the indegree distributions for the TUCs and the random collections.

(2) *Outdegree distributions*: Outdegree distributions seem to not follow the power law at small values. However, larger values do seem to follow such a distribution, resulting in a “drooping head” of the log-log plot as observed in earlier work. A good characterization of outdegrees for the web graph has not yet been offered, especially one that would satisfactorily explain the drooping head.

(3) *Connected component sizes*: (cf. Section 2) We consider the size of the largest strongly-connected component, the second-largest, third-largest and so forth as a distribution, for each graph of interest. We consider similar statistics for the sizes of weakly-connected components. Specifically, we will show that they obey power laws at all scales, and study the exponents of the power law (called *SCC/WCC exponent*). We also report the ratio of the size of the largest strongly-connected component to the size of the largest weakly-connected component. For the significance of these parameters, we refer the reader to [13], and note that the location of a web page in the connected component decomposition crucially determines the reachability of this page (often related to its popularity).

(4) *Bipartite cores*: Bipartite cores are graph-theoretic signatures of community structure on the web. A $K_{i,j}$ *bipartite core* is a set of $i + j$ pages such that each of i pages contains a hyperlink to all of the remaining j pages. We pick representative values of i and j , and focus on $K_{5,7}$'s, which are sets of 5 “fan” nodes, each of which points to the same set of 7 “center” nodes. Since computing the exact number of $K_{5,7}$'s is a complex subgraph enumeration problem that is intractable using known techniques, we instead estimate the number of node-disjoint $K_{5,7}$'s for each graph of interest. To perform this estimation, we use the techniques of [28, 29]. The number of communities (cores) is an estimate of community structure with the TUC. The $K_{5,7}$ *factor* is the relative size of the community to the size of the nodes that participate in $K_{5,7}$'s in it. The higher the factor, the less one can view the TUC as a single well defined community.

(5) *URL compressibility and namespace utilization*: The URL namespace can be viewed as a tree, with the root node being represented by the null string. Each node of the tree corresponds to a URL prefix (say `www.foo.com`) with all URLs that share that prefix, (e.g, `www.foo.com/bar` and `www.foo.com/rab`) being in the subtree subtended at that node. For each subgraph and each value d of the depth, we study the following distribution: for each s , the number of depth- d nodes whose subtrees have s nodes. We will see that these follow the power law. Following conventional source coding theory, it follows that this skew in the population distributions of the URL namespace can be used to design improved compression algorithms for URLs. The details of this analysis are beyond the scope of the present paper.

3.3 Experimental infrastructure

We performed these experiments on a small cluster of Linux machines with about 1TB of disk space. We created a number of data sets from two original sets of pages. The first set consists of about 500K pages from the IBM intranet. We treat this data as a single entity, mainly for purposes of comparison with the external web. The second set consists of 60M pages from the web at large, crawled in Oct. 2000. These 60M pages represent approximately 750GB of content. The crawling algorithm obeyed all politeness rules, crawling no site more often than once per second. Therefore, while we had collected 750GB of content (crawling about 1.3M sites) no more than 12K pages had been crawled from any one site.

4 Results and interpretation

Our results are shown in the following tables and figures. Though we have an enormous amount of data, we try to present as little as possible, while conveying the main thoughts. All the graphs here refer to node-induced subgraphs and the arcs refer to the arcs in the induced subgraph. Our tables show the parameters in terms of the graphs while our figures show the consistency of the parameters across different graphs, indicating a fractal nature.

Table 1 shows all the parameters for the STREAM1 through STREAM7. The additional parameter, *expansion factor*, refers to the fraction of hyperlinks that point to nodes in the same collection to the total number of hyperlinks. As we can see, the numbers are quite consistent with earlier work. For instance, the indegree exponent is -2.1, the SCC exponent is around -2.15, and the WCC exponent is around -2.3. As we can see, the ratios of IN, OUT, SCC with respect to WCC are also consistent with earlier work.

Table 2 shows the results for the three special graphs: INTRANET, HOSTGRAPH, and GEO. The expansion factor for the INTRANET is 2.158 while the indegree exponent is very different from that of other graphs. The WCC exponent for HOSTGRAPH is not meaningful since there is a single component that is 99.4% of the entire graph.

Table 3 shows the results for single keyword queries. The graphs in the category are in few hundreds of thousands. Table 4 shows the results for double keyword graphs. The graphs in this category are in few tens of thousands. A specific interesting case is the large $K_{5,7}$ factor for the keyword MATH, which probably arises since pages containing the term MATH is probably not a TUC since it is far too general.

Table 5 shows the averaged results for the 100 sites SUBDOMAIN1, ..., SUBDOMAIN100.

Next, we point out the consistency of the parameters across various graphs. For ease of presentation, we picked a small set of TUCs and plotted the distribution of indegree, outdegree, SCC, WCC on a log-log scale (see Figures in Appendix). Figure 2 shows the indegree and outdegree distributions for five of the TUCs.

As we see, the shape of plots are strikingly alike. As observed in earlier studies, a drooping initial segment is observed in the case of outdegree. Figure 3 shows the component distributions for the graphs. Again, the similarity of shapes is striking. Figure 4 show the URL tree size distribution. The figures show remarkable self-similarity that exists both across graphs and within each graph across different depths.

4.1 Discussion

We now mention four interesting observations based on the experimental results. Following [13] (see also Section 2), we say that a slice of the web graph *has the bow-tie structure* if the SCC, IN, and OUT, each accounts for a large constant fraction of the nodes in the slice.

(1) Almost all nodes (82%) of the HOSTGRAPH are contained in a giant SCC (Table 2). This is not surprising, since one would expect most websites to have at least one page that belongs to the SCC.

(2) The (microscopic) local graphs of SUBDOMAIN1, ..., SUBDOMAIN100, look surprisingly like the web graph (see Table 5. Each has an SCC flanked by IN and OUT sets that, for the most part, have sizes proportional to their size on the web as a whole, about 40% for the SCC, for instance. Large websites seemed to have a more clearly defined bow-tie structure than the smaller, less developed ones.

(3) Keyword based TUCs corresponding to KEYWORD1, ..., KEYWORD10 (see Tables 3 and 4) exhibit similar phenomena; the differences often owe to the extent to which a community has a well-established presence on the web. For example, it appears from our results that the GOLF is a well-established web community, while RESTAURANT is a newer developing community on the web. While the mathematics community had a clearly defined bow-tie structure, the less developed geometry community lacked one.

(4) Considering STREAM1, ..., STREAM7, we find the surprising fact (Table 1) that the union of a random collection of TUCs contains a large SCC. This shows that the SCC of the web is very resilient to node deletion and does not depend on the existence of large taxonomies (such as yahoo.com) for its connectivity. Indeed, as we remarked earlier, each of these streams contain very few arcs which are not entirely local to the website. However, the bow-tie structure of each website allows the few intersite arcs to be far more valuable than one would expect.

4.2 Analysis and summary

The foregoing observation about the SCC of the streams, while surprising, is actually a direct consequence of the following theorem about random edges in graphs with large strongly connected components.

Theorem 1. *Consider the union of n/k graphs on k nodes each, where each graph has a strongly connected component of size ek . Suppose we add dn arcs whose heads and tails are uniformly distributed among the n nodes, then provided that d is at least of the order $1/(ek)$, with high probability, we will have a strongly connected component of size of the order of en on the n -node union of the n/k graphs.*

The proof of Theorem 1 is fairly straightforward. On the web, n is about 1 billion, k , the size of each TUC, is about 1 million (in reality, there are more than 1K TUCs that overlap, which only makes the connectivity stronger), and e is about $1/4$. Theorem 1 suggests that the addition of a mere few thousand arcs scattered uniformly throughout the billion nodes will result in very strong connectivity properties of the web graph!

Indeed, the evolving copying models for the web graph proposed in [27] incorporates a uniformly random component together with a copying stochastic process. Our observation above is, in fact, lends consid-

Nodes $\times 10^6$	Arcs $\times 10^6$	Expansion factor	Indeg. exp.	Outdeg. exp.	SCC exp.	WCC exp.	WCC $\times 10^6$	SCC/ WCC	IN/ WCC	OUT/ WCC	$K_{5,7}$ factor
6.55	46.8	2.06	-2.07	-2.12	-2.16	-2.32	4.69	0.24	0.23	0.23	47.2
6.47	45.7	2.06	-2.08	-2.24	-2.14	-2.28	4.60	0.23	0.19	0.24	50.1
6.38	48.1	2.05	-2.06	-2.15	-2.15	-2.24	4.47	0.24	0.20	0.23	49.5
6.84	50.0	2.04	-2.12	-2.30	-2.14	-2.27	4.86	0.23	0.21	0.23	43.5
6.83	48.2	2.06	-2.08	-2.27	-2.11	-2.29	4.90	0.24	0.20	0.23	45.4
6.77	49.3	2.01	-2.10	-2.32	-2.11	-2.25	4.78	0.23	0.20	0.24	45.3
6.23	43.5	2.03	-2.13	-2.19	-2.15	-2.27	4.31	0.22	0.19	0.23	46.9

Table 1: Results for STREAM1 through STREAM7.

Subgraph	Nodes $\times 10^3$	Arcs $\times 10^3$	Indeg. exp.	SCC exp.	WCC exp.	WCC $\times 10^3$	SCC/ WCC	IN/ WCC	OUT/ WCC	$K_{5,7}$ factor
INTRANET	285.5	1910.7	-2.31	-2.53	-2.83	207.7	0.20	0.48	0.17	56.13
HOSTGRAPH	663.7	1127.9	-2.34	-2.81		659.9	0.82	0.04	0.13	72.64
GEO	410.7	1477.9	-2.51	-2.69	-2.27	2.1	0.87	0.03	0.10	139.9

Table 2: Results for graphs: INTRANET, HOSTGRAPH, and GEO.

Subgraph	Nodes $\times 10^3$	Arcs $\times 10^3$	Indeg. exp.	SCC exp.	WCC exp.	WCC $\times 10^3$	SCC/ WCC	$K_{5,7}$ factor
BASEBALL	336.5	3444.4	-2.09	-2.16	-2.30	33.2	0.12	55.85
GOLF	696.8	8512.8	-2.06	-2.06	-2.18	47.3	0.15	44.48
MATH	831.7	3787.8	-2.85	-2.66	-2.73	50.2	0.28	148.7
MP3	497.3	7233.2	-2.20	-2.39	-2.20	47.6	0.28	57.18
RESTAURANT	623.0	3592.5	-2.33	-2.47	-2.28	7.96	0.31	115.2

Table 3: Results for single keyword query graphs KEYWORD1 through KEYWORD5.

Subgraph	Nodes $\times 10^3$	Arcs $\times 10^3$	WCC $\times 10^3$	SCC/ WCC	$K_{5,7}$ factor
BASEBALL YANKEES	24.0	320.0	3813	0.73	45.82
GOLF TIGER WOODS	14.9	62.8	1501	0.20	83.02
MATH GEOMETRY	44.0	86.9	1903	0.27	407.52
MP3 NAPSTER	27.1	321.4	1775	0.36	35.19
RESTAURANT SUSHI	7.4	23.7	167	0.72	132.14

Table 4: Results for double keyword query graphs KEYWORD6 through KEYWORD10.

Nodes $\times 10^3$	Arcs $\times 10^3$	WCC $\times 10^3$	SCC/ WCC	$K_{5,7}$ factor
7.17	108.42	7.08	0.42	22.97

Table 5: Averaged results for SUBDOMAIN1 through SUBDOMAIN100.

erable support to the legitimacy of this model. These observations, together with Theorem 1, imply a very interesting detailed structure for the SCC of the webgraph.

The web comprises several thematically unified clusters (TUCs). The common theme within a TUC is one of many diverse possibilities. Each TUC has a bow-tie structure that consists of a large strongly connected component (SCC). The SCCs corresponding to the TUCs are integrated, via the navigational backbone, into a global SCC for the entire web. The extent to which each TUC exhibits the bow-tie structure and the extent to which its SCC is integrated into the web as a whole indicate how well-established the corresponding community is.

An illustration of this characterization of the web is shown in Figure 1.

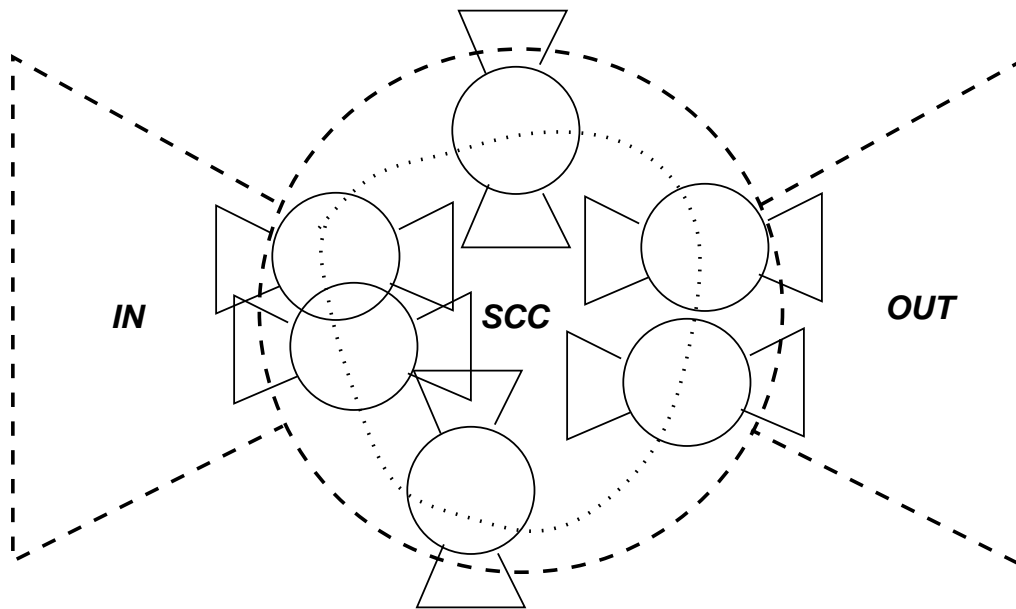


Figure 1: TUCs connected by the navigational backbone inside the SCC of the web graph.

5 Conclusions

In this paper, we have examined the structure of the web in greater detail than earlier efforts. The primary contribution is two-fold. First, the web exhibits self-similarity in several senses, at several scales. The self-similarity is pervasive, in that it holds for a number of parameters. It is also robust, in that it holds irrespective of which particular method is used to carve out small subgraphs of the web. Second, these smaller thematically unified subgraphs are organized into the web graph in an interesting manner. In particular, the local strongly connected components are integrated into the global SCC. The connectivity of the global SCC is very resilient to random and large scale deletion of websites. This indicates a great degree of fault-tolerance on the web, in that there are several alternate paths between nodes in the SCC.

While our understanding of the web as a graph is greater now than ever before, there are many lacunae in our current understanding of the graph-theoretic structure of the web. One of the principal holes deals

with developing stochastic models for the evolution of the web graph (extending [27]) that are rich enough to explain the fractal behavior of the web in such amazingly diverse ways and contexts.

Acknowledgments

Thanks to Raymie Stata and Janet Wiener (Compaq SRC) for some of the code. The second author thanks Xin Guo for her encouragement to this project.

References

- [1] W. Aiello, F. Chung, and L. Lu. A random graph model for massive graphs. *Proc. 32d STOC*, 1999.
- [2] L. Adamic and B. Huberman. The nature of markets on the world wide web. *Xerox PARC Technical Report*, 1999.
- [3] L. Adamic and B. Huberman. Scaling behavior on the world wide web. *Technical comment on Barabasi and Albert 99*.
- [4] M. Adler and M. Mitzenmacher. Towards compressing web graphs. *Proc. IEEE Data Compression Conference*, 2001, To appear.
- [5] A. Arasu, A. Tomkins, and J. Tomlin. Pagerank Computation and the Structure of the Wweb: Experiments and Algorithms. Submitted, 2001.
- [6] S. Abiteboul, D. Quass, J. McHugh, J. Widom, and J. Wiener. The Lorel Query Language for Semistructured Data. *Int. J. on Digit. Libr.*, 1(1):68–88, 1997.
- [7] A. Barabasi and R. Albert. Emergence of scaling in random networks. *Science*, 286(509), 1999.
- [8] G. Arocena, A. Mendelzon, and G. Mihaila Applications of a Web Query Language In *Proc. 6th Int'l. WWW Conf.*, Santa Clara, April 1997.
- [9] P. Barford, A. Bestavros, A. Bradley, and M. E. Crovella. Changes in web client access patterns: Characteristics and caching implications. in *World Wide Web, Special Issue on Characterization and Performance Evaluation*, 2:15-28, 1999.
- [10] K. Bharat and M. Henzinger. Improved algorithms for topic distillation in hyperlinked environments. *Proc. 21st SIGIR*, 1998.
- [11] B. Bollobas. *Random Graphs*. Academic Press, 1985.
- [12] S. Brin and L. Page. The anatomy of a large scale hypertextual web search engine. *Proc. 7th WWW*, 1998.
- [13] A. Broder, R. Kumar, F. Maghoul, P. Raghavan, S. Rajagopalan, R. Stata, A. Tomkins, and J. Wiener. Graph Structure in the web. *Proc. 9th WWW*, 2000.
- [14] R. A. Botafogo and B. Shneiderman. Identifying aggregates in hypertext structures. *Proc. 3rd ACM Conference on Hypertext*, 1991.
- [15] J. Carriere and R. Kazman. webQuery: Searching and visualizing the Web through connectivity. *Proc. 6th WWW*, 1997.
- [16] S. Chakrabarti, B. Dom, D. Gibson, J. Kleinberg, P. Raghavan, and S. Rajagopalan. Automatic resource compilation by analyzing hyperlink structure and associated text. *Proc. 7th WWW*, 1998.

- [17] S. Chakrabarti, B. Dom, D. Gibson, S. Ravi Kumar, P. Raghavan, S. Rajagopalan, and A. Tomkins. Experiments in topic distillation. *Proc. ACM SIGIR workshop on Hypertext Information Retrieval on the Web*, 1998.
- [18] S. Chakrabarti, M. van den Berg, and B. Dom. Focused crawling: A new approach to topic-specific web resource discovery. *Proc. 8th WWW*, 1999.
- [19] S. Chakrabarti, D. Gibson, and K. McCurley. Surfing the web backwards. *Proc. 8th WWW*, 1999.
- [20] S. Deerwester, S. Dumais, G. Furnas, T. Landauer, and R. Harshman. Indexing by latent semantic analysis. *J. ASIS*, 41(6), 391–407, 1990.
- [21] M. Faloutsos, P. Faloutsos, and C. Faloutsos. On power law relationships of the internet topology. *Proc. ACM SIGCOMM*, 1999.
- [22] S. Glassman. A caching relay for the world wide web. *Proc. 1st WWW*, 1994.
- [23] F. Harary. *Graph Theory*. Addison Wesley, 1975.
- [24] B. Hill. A Simple method for inferring the tail behavior of distributions. *Annals of Statistics*, 1975.
- [25] B. Huberman, P. Pirollo, J. Pitkow, and R. Lukose. Strong regularities in world wide web surfing. *Science*, 280:95-97, 1998.
- [26] J. Kleinberg. Authoritative sources in a hyperlinked environment. *J. ACM*, 2000.
- [27] R. Kumar, P. Raghavan, S. Rajagopalan, D. Sivakumar, A. Tomkins, and E. Upfal. Random graph models for the web graph. *Proc. 41st FOCS*, 2000.
- [28] R. Kumar, P. Raghavan, S. Rajagopalan, and A. Tomkins. Trawling the web for cyber communities. *Proc 8th WWW*, 1999.
- [29] R. Kumar, P. Raghavan, S. Rajagopalan, and A. Tomkins. Extracting large scale knowledge bases from the web. *Proc. VLDB*, 1999.
- [30] R. Kumar, P. Raghavan, S. Rajagopalan, and A. Tomkins. Human effort in semi-automated taxonomy construction. *IBM Research Report*, 2000.
- [31] <http://linkage.rockefeller.edu/wli/zipf/>.
- [32] R. M. Lukose and B. Huberman. Surfing as a real option. *Proc. 1st Intl. Conf. Information and Computation Economies*, 1998.
- [33] C. Martindale and A. K. Konopka. Oligonucleotide frequencies in DNA follow a Yule distribution. *Computer & Chemistry*, 20(1):35-38, 1996.
- [34] A. Mendelzon, G. Mihaila, and T. Milo. Querying the world wide web. *J. of Digital Libraries*, 1(1), pp. 68-88, 1997.
- [35] A. Mendelzon and P. Wood. Finding regular simple paths in graph databases. *SIAM J. Comp.*, 24(6):1235-1258, 1995.
- [36] C. Papadimitriou, P. Raghavan, H. Tamaki, and S. Vempala. Latent Semantic Indexing: A Probabilistic Analysis. *Proc. 17th PODS*, 1998.
- [37] E. M. Palmer. *Graphical Evolution*. John Wiley, 1985.
- [38] V. Pareto. *Cours d'economie politique*. Rouge, Lausanne et Paris, 1897.

- [39] P. Pirolli, J. Pitkow, and R. Rao. Silk from a sow's ear: Extracting usable structures from the web. *Proc. ACM SIGCHI*, 1996.
- [40] J. Pitkow and P. Pirolli. Life, death, and lawfulness on the electronic frontier. *Proc. ACM SIGCHI*, 1997.
- [41] H. A. Simon. On a class of skew distribution functions. *Biometrika*, 42:425-440, 1955.
- [42] E. Spertus and L. Stein. A Hyperlink-Based Recommender System Written in Squeal. Workshop on Web Information and Data Management, November 6, 1998.
- [43] E. Spertus. ParaSite: Mining Structural Information on the web. *Proc. 6th WWW*, 1997.
- [44] H. D. White and K. W. McCain. Bibliometrics. In *Ann. Rev. Info. Sci. and Technology*, 119–186, 1989.
- [45] G. U. Yule. *Statistical Study of Literary Vocabulary*. Cambridge University Press, 1944.
- [46] G. K. Zipf. *Human Behavior and the Principle of Least Effort*. Addison-Wesley, 1949.