

Efficient Parallel Computation of PageRank

Christian Kohlschütter, Paul-Alexandru Chirita, and Wolfgang Nejdl

L3S Research Center / University of Hanover
Deutscher Pavillon, Expo Plaza 1
30539 Hanover, Germany
{kohlschuetter,chirita,nejdl}@l3s.de

Abstract. PageRank inherently is massively parallelizable and distributable, as a result of web’s strict host-based link locality. In this paper we show that the Gauß-Seidel iterative method for solving linear systems can be successfully applied in such a parallel ranking scenario in order to improve convergence. By introducing a two-dimensional web model and by adapting the PageRank to this environment, we present and evaluate efficient methods to compute the exact rank vector even for large-scale web graphs in only a few minutes and iteration steps, with intrinsic support for incremental web crawling, and without the need for page sorting/reordering or for sharing global information.

1 Introduction

Search engines are the enabling technology for finding information on the Internet. They provide regularly updated snapshots of the Web and maintain a searchable index over all retrieved web pages. Its size is currently reaching several billions of pages from about 54 million publicly accessible hosts, but these amounts are rapidly increasing [17]. When searching such huge datasets, one would usually receive quite a few pages in response to her query, some of them being much more relevant than others. This gave birth to a lot of ordering (or ranking) research, the most popular algorithm being Google’s PageRank [18], which recursively determines the importance of a web page by the importance of all the pages pointing to it.

Although improvements for a centralized computation of PageRank have been researched in detail [1, 9, 11, 5, 12, 20, 15], approaches on distributing it over several computers have caught researchers’ attention only recently. In this paper we introduce a new approach to computing the exact PageRank in a parallel fashion. We obtain exact results faster than all the other existing algorithms, improving by several orders of magnitude over the other algorithms generating exact PageRank scores. We achieve this by modeling the web graph in a two-dimensional fashion (with the URL’s hostname as the primary criterion), thus separating it into reasonably disjunct partitions, which are then used for distributed, incremental web crawling [6] and PageRank computation.

The remainder of the paper is organized as follows. After reviewing the PageRank algorithm, common web graph representation techniques and existing parallel versions of PageRank in Section 2, we introduce our two-dimensional

web graph model in Section 3. We then present a refined PageRank algorithm in Section 4, and show that the convergence improvements of the Gauß-Seidel method for solving linear systems can also be efficiently applied in a parallelized PageRank scenario. Experimental results are discussed in Section 5. Finally, Section 6 concludes with discussion of further work.

2 Background and Previous Work

2.1 Web Graph Representation

Computing the PageRank vector using a materialized in-memory adjacency matrix for a large web graph is definitely not feasible. A common solution is to store the links in a format like “Destination Page ID, Out-degree, Source Page IDs...” (which resembles L). Because pages only link to a few others (the link matrix is sparse), this results in much lower memory requirements of the link structure, in the magnitude of $|L| \cdot \bar{n}^{-1}$ (\bar{n} = average outdegree). Of course, compression techniques [14] or disk-based “swapping” [9, 5] can improve the space requirements even further. But with the permanent growth of the web, even such techniques will soon hit memory limits of a single computer, or unacceptably slow down the computation process. In this paper, we thus propose a new storage format for the adjacency matrix of the web graph, based on the separation between global (host) and local information about each page.

2.2 PageRank

The main concept behind the PageRank paradigm [18] is the propagation of importance from one Web page towards others, via its out-going (hyper-)links. Each page $p \in P$ (P is the set of all considered pages) has an associated rank score $r(p)$, forming the rank vector \mathbf{r} . Let L be the set of links, where (s, t) is contained iff page s points to page t and $L(p)$ be the set of pages p points to (p ’s outgoing links). The following iteration step is then repeated until all scores r stabilize to a certain degree:

$$\forall t \in P : r^{(i)}(t) = (1 - \alpha) \cdot \tau(t) + \alpha \sum_{(s,t) \in L} \frac{r^{(i-1)}(s)}{|L(s)|} \quad (1)$$

The formula consists of two portions, the jump component (left side) and the walk component (right side), weighted by α (usually 0.85). $r^{(i-1)}(s) \cdot |L(s)|^{-1}$ is the uniformly distributed fraction of importance a page s can offer to one of its linked pages t for iteration i . Intuitively, a “random surfer” will follow an outgoing link from the current page (walk) with probability α and will get bored and select a random page (jump) with probability $(1 - \alpha)$. The main utility of α is however to guarantee convergence and avoid “rank sinks” [3].

This “random-walk” is in fact the interpretation of the Markov chain associated to the web graph, having \mathbf{r} as the state vector and A (see Equation 2) the transition probability from one page to another. We can therefore also write Equation 1 in matrix terms as follows:

$$\mathbf{r} = (1 - \alpha) \cdot \boldsymbol{\tau} + \alpha A \mathbf{r} \quad (2)$$

Equation 1 also represents the linear system representation of this computation using the Jacobi method. This enables the consideration of using other stationary iterative solvers, such as the Gauß-Seidel method, which was said not to be efficiently parallelizable here [1, 5]. Actually, there already are parallel Gauss-Seidel implementations for certain scenarios such as the one described in [13], using block-diagonally-bordered matrices; however, they all admit their approach was designed for a static matrix; after each modification, a specific preprocessing (sorting) step is required, which can take longer than the real computation. Because the web is highly dynamic, almost 40% of all links change in less than one week [6]. Thus, disregarding this preparation step veils the real overall processing time. Steady reorganization of coordinates in a huge link matrix simply imposes an unjustified management overhead.

2.3 Other Parallel PageRank Algorithms

Existing approaches to PageRank parallelization can be divided into two classes: Exact Computations and Approximations.

Parallel Computations. In this scenario, the web graph is initially partitioned into blocks: grouped randomly (e.g., P2P PageRank [19]), lexicographically sorted by page (MIKELab PageRank [16] and Open System PageRank [21]) or balanced according to the number of links (PETSc PageRank [8]). Then, standard iterative methods such as Jacobi (Equation 1) or Krylov subspace [8] are performed over these pieces in parallel. The partitions periodically must exchange information: Depending on the strategy this can expose suboptimal convergence speed because of the Jacobi method and result in heavy inter-partition I/O (e.g., in MIKELab PageRank, computing the rank for a page t requires access to all associated source page ranks $r(s)$ across all partitions).

PageRank Approximations. The main idea behind these approaches is that it might be sufficient to get a rank vector which is comparable, but not equal to PageRank. Instead of ranking pages, higher-level formations are used, such as the inter-connection/linkage between hosts, domains, server network addresses or directories, which is orders of magnitudes faster. The inner structure of these formations (at page level) can then be computed in an independently parallel manner (“off-line”), as in BlockRank [10], SiteRank [24], the U-Model [4], ServerRank [23] or HostRank/DirRank [7].

In our approach, we will try to take the best out of both approaches: the exactness of a straight PageRank computation but the speed of an approximation, without any centralized re-ranking.

3 The Two-Dimensional Web

3.1 Host-based Link Locality

Bharat et al. [2] have shown that there are two different types of web links dominating the web structure, “intra-site” links and “inter-site” ones. A “site” can be a domain (.yahoo.com), a host (geocities.yahoo.com) or a directory

on a web server (<http://www.geocities.com/someuser/>). In general, we can define a site as an interlinked collection of pages identified by a common name (domain, host, directory etc.), and under the control of the same authority (an authority may of course own several sites).

Due to web sites' hypertext-navigable nature, it is supposable that a site contains more internal than external links. In fact, about 93.6% of all non-dangling links are intra-host and 95.2% intra-domain [10]. This assumed block structure has been visualized by Kamvar et al. [10] using dotplots of small parts (domain-level) of the "LargeWeb" graph's link matrix [22]. In these plots, the point (i, j) is black, if there is a link from page p_i to p_j , clear otherwise.

We performed such a plot under the same setting, but on whole-graph scale. The outcome is interesting: a clear top-level-domain (TLD) dominant structure (see Figure 1a). For example, the .com TLD represents almost 40% of the complete structure and has high connectivity with .net and .org, whereas the .jp domain shows almost no interlinkage with other TLDs. However, if we only inspect the .com domain (see Figure 1b, the dotplot depicts a diagonally dominant structure. The diagonal represents links from target pages near by the source page (which are *inter-host* pages). Both results are primarily caused by the lexicographical order of URLs, where hostnames are reversed.

But is this costly *sorting* over all URLs necessary at all? To further analyze the impact of hostname-induced link locality, we redraw the LargeWeb dotplot in a *normalized* (histographical) fashion, where a dot's greyscale value depicts the cumulative percentage of links in a specific raster cell. In addition, we do not sort the pages lexicographically, but only group them per host and permute all hosts randomly to exclude any lexicographical relationship between them. The clear diagonal dominance now also becomes visible on whole-graph scale (Figure 1c).

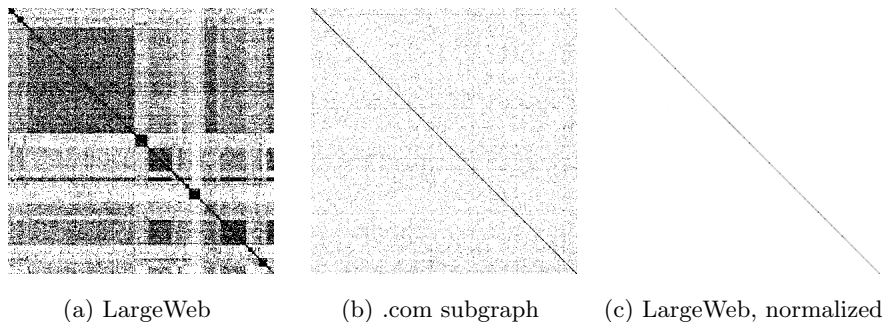


Fig. 1. Linkage dotplots, sorted by URL.

3.2 From Numbers to Tuples

It should be obvious that the web was already designed to be two-dimensional: Hostnames are “namespaces” aimed to disambiguate different local contexts (i.e., paths like “/dir/ index.html”). Previous approaches to web graph partitioning always resulted in having *one* unique ID associated to each page. Such a single page ID provides a very compact representation of the web graph, which can be visualized in a matrix dotplot as shown above. But it also requires continuous reorganization (resorting) for newly added pages. Otherwise, a mixture of hosts along the URL IDs would render a host no longer characterizable by a closed interval of IDs, thereby losing the advantage of link locality. One may introduce gaps in the numbering to reduce the sorting costs, but still, all subsequent pages will have to be renumbered once the gap is filled. In a distributed scenario, this can cause extensive network I/O by repeatedly moving pages from one partition to another.

We therefore propose a different page identification scheme, based on the affiliation of each page to a specific host and independently of pages from other hosts. More specifically, we propose using a *tuple* consisting of two independent, positive integers, a HostID (only dependent on the URL’s hostname) and a LocalID (only identifying the remaining local components – path and query string). The addition of new local pages to a specific host, as well as of new hosts, is very easy, since renumbering is no longer necessary.

As an implementation-specific note, we expect that for current web graphs, it is sufficient to store the tuples as two `uint32` four-byte integers. We then can address a maximum of 4.29 billion hosts and a maximum of 4.29 billion pages per host in 8 bytes. For small hosts, we could even reduce the local part to 16 bit, thereby further cutting down memory footprint.

4 Partitioned PageRank

We will now consider the impact of such a partitioning scheme on the PageRank algorithm. We will first present an analysis that unifies two of the most common algorithms for solving linear systems, Gauß-Seidel and Jacobi. Then, we will apply this analysis to propose an improved parallel PageRank algorithm, and finally we will discuss several optimization issues.

4.1 Unifying Jacobi and Gauss-Seidel

It has been observed that the Gauß-Seidel iteration method compared to the Jacobi method can speed-up PageRank convergence by a factor of 2, as it uses scores of the current iteration as soon as they become available [1]:

$$\forall (s, t) \in L: \quad r^{(i)}(t) = (1 - \alpha) \tau(t) + \alpha \left(\sum_{s < t} \frac{r^{(i)}(s)}{|L(s)|} + \sum_{s > t} \frac{r^{(i-1)}(s)}{|L(s)|} \right) \quad (3)$$

As opposed to the Jacobi iteration, the Gauß-Seidel variant requires iterating over the links $(s, t) \in L$ in a strictly ascending order. At first glance, this seems to be a major drawback when we want to apply it to a distributed, partitioned web graph. To clarify the impact of the restriction of link order, we derive a common base algorithm for both, Jacobi (equation 1) and Gauß-Seidel (equation 3) algorithms: We define an intermediate ranking vector $r^{(i-1,i)}$ that combines the vectors of the previous and the current iteration, depending on the state of a ranked page p in the set of available pages P ($P = P' \cup P''$; $\nexists p : p \in P' \wedge p \in P''$; P' contains all pages which have already been ranked for iteration i ; P'' contains all other pages, whose score has not been touched since iteration $i - 1$):

$$r^{(i-1,i)}(p) := \begin{cases} r^{(i)}(p) & \text{if } p \in P' \\ r^{(i-1)}(p) & \text{if } p \in P'' \end{cases}; r^{(i)}(t) = (1 - \alpha) \tau(t) + \alpha \sum_{(s,t) \in L} \frac{r^{(i-1,i)}(s)}{|L(s)|} \quad (4)$$

Under this setting, for the Gauß-Seidel method, $P' = \{ p \mid p < k \}$ and $P'' = \{ p \mid p \geq k \}$, with $k \in \{1, 2, \dots, |P|\}$, whereas for the Jacobi method, we have $P' = \emptyset$ and $P'' = P$. Both iteration methods, Jacobi and Gauß-Seidel, can then be simplified to this joint formula:

$$r^{(*)}(t) = (1 - \alpha) \tau(t) + \alpha \sum_{(s,t) \in L} \frac{r^{(*)}(s)}{|L(s)|}, \text{ with } r^{(*)}(t) = r^{(i-1,i)}(t) \quad (5)$$

From Equation 4, we know that before each iteration i , $\mathbf{r}^{(*)} = \mathbf{r}^{(i-1)}$ and after the iteration $\mathbf{r}^{(*)} = \mathbf{r}^{(i)}$. The state of $\mathbf{r}^{(*)}$ during the iteration then only depends on the order of links $(s, t) \in L$ (the way how P' and P'' are determined). This iteration method has worst-case convergence properties of Jacobi and best-case of Gauß-Seidel, depending on the order of elements, random order vs. strictly ascending order, while always providing the same per-iteration running time as the Jacobi iteration.

We further generalize the impact of the rules for P' and P'' : We argue that if only a small fraction F of all links concerned ($|F| \ll |L|$) is not in strictly ascending order, the overall convergence speed still remains in the magnitude of standard Gauß-Seidel. In our case, in order to be able to parallelize the Gauß-Seidel algorithm, we will assign inter-host/inter-partition links (about 6%) to this small fraction.

4.2 Reformulating PageRank

For such an optimization, let us reformulate our above mentioned unified PageRank equation using our new two-dimensional page numbering scheme. Thus, page variables “ p ” will be replaced by page tuples “ (p_x, p_y) ”, with p_x representing the page’s HostID and p_y its LocalID. Also, to account for the separation of inter- and intra-host links, the formula becomes:

$$\begin{aligned}
r^{(*)}(\mathbf{t}) &= (1 - \alpha) \tau(\mathbf{t}) + \alpha \left(v_I^{(*)}(\mathbf{t}) + v_E^{(*)}(\mathbf{t}) \right) \\
v_I^{(*)}(\mathbf{t}) &= \sum_{(s,\mathbf{t}) \in L} \frac{r^{(*)}(\mathbf{s})}{|L(s)|} \quad \forall \text{host}(s) = \text{host}(t) \\
v_E^{(*)}(\mathbf{t}) &= \sum_{(s,\mathbf{t}) \in L} \frac{r^{(*)}(\mathbf{s})}{|L(s)|} \quad \forall \text{host}(s) \neq \text{host}(t)
\end{aligned} \tag{6}$$

Since $v_I^{(*)}(\mathbf{t})$ solely requires access to *local* rank portions of $\text{host}(t)$, it can efficiently be computed from scores stored in RAM. The local problem of ranking intra-host pages is solvable via a fast, non-parallel Gauß-Seidel iteration process. There is no need for intra-host vote parallelization – instead, we parallelize on the host-level, thus necessitating only inter-host communication, which is limited to the exchange of external votes.

Our approach produces the same ranks as the original PageRank, while being more scalable than the other parallel PageRank algorithms. This is mainly due to the parallelization of the Gauß-Seidel algorithm, in which we take advantage of web’s host-oriented block structure.

4.3 Reaching Optimal Performance

Optimizing Communication Cost. While votes between hosts of the same partition (server) can easily be conveyed in RAM, votes between hosts of different partitions require network communication. The gross total for exchanging external votes over the network must not be underestimated. In our setup with the LargeWeb graph, almost 33 million votes must be exchanged between partitions. For bigger web graphs, this could rise up to a few billion and can easily lead to network congestion if too much information is transmitted per vote.

As opposed to other approaches, where a vote consisted of target page ID (sometimes along with source page ID) and score, we simply reduced this to transmitting a single score value per page, because in a static graph, the link structure itself does not change during the iteration cycle. More generally, the link structure of all the pages that exchange votes between two partitions pages only needs to be determined whenever the graph changes (in the case of incremental web crawling) and then to be sent to the specific target partition. Moreover, the source page does not need to be specified in order to compute the PageRank score (see Equation 6), but only the target page ID. Additionally, by grouping the list of target pages by host, we need to transmit each HostID only once.

Most notably, each partition has to transmit only one single value per target page, not per link to thatpage, since all votes from local pages that link to a specific page can be aggregated to a single value (surprisingly, this simple, but very effective approach did not appear in any previous work):

$$v_E^{(*)}(\mathbf{t}) = \sum_{\beta \in \Pi} \sum_{(s,\mathbf{t}) \in L_\beta} \frac{r^{(*)}(\mathbf{s})}{|L|} = \sum_{\beta \in \Pi} v_\beta^{(*)}(\mathbf{t}) \quad \forall \text{host}(s) \neq \text{host}(t) \tag{7}$$

with Π being the set of partitions containing links towards t , and β each of these partitions.

Transferring $v_\beta(t)$ (the sum of votes from partition L_β to t) as a single value will then reduce the network load dramatically. Using this optimization, we can show a reduction of vote exchanges by 89% with the DNR-LargeWeb graph. Table 1 depicts the difference between inter-partition links and votes and their quota of all links.

Type	Amount	Percent
Total Links	601,183,777	100%
Inter-Partition Links	32,716,628	5.44%
Inter-Partition Votes	3,618,335	0.6%

Table 1. LargeWeb Inter-Partition links and votes

Computational Load Balancing. In order to keep the convergence behavior of the centralized PageRank in our parallel scenario, inter-partition votes must be exchanged after every iteration (see [16] for a discussion of consequences of not doing so). To keep the overall computation time still low, all intra-partition computations and after that all network communication should terminate isochronously (at the same time). Because intra-partition computation is directly proportional to the number of pages per partition (see Equation 6), this either means that all available servers must be equally fast, or the graph has to be at least partitioned adequately to the performance of the servers. Moreover, other slow-down factors could also influence the running time, such as different network throughput rates of cheap NICs and system boards (even with the same nominal speed).

A good strategy to load-balancing Parallel PageRank in a heterogeneous environment could be running a small test graph on all new servers, measure computation speeds, and balance the real graph accordingly. In any case, memory overflows due to bad balancing parameters like in PETSc PageRank are avoided, and no manual interaction to find these parameters is necessary.

5 Experiments

We first converted the Stanford DNR-LargeWeb graph [22] into the new tuple representation, resulting in 62.8M pages and 601M links distributed over 470,000 hosts with averaged 137.5 pages each (maximum was 5084 pages per host); the inter-host link percentage¹ is 6.19% (see Table 2).

For our PageRank experiments, we sorted sorted the available hosts by their page count in descending order and distributed the pages host-wise in a round-robin manner over 8 partitions of equal size ($\frac{1}{8}$ of the graph just fitted into our smallest server’s RAM).

¹ Unfortunately, the last 8 million pages of DNR-LargeWeb could not be converted, since there was no URL associated with them – thus, our numbers slightly differ from the ones in [10].

Although the pages-per-host distribution was not strictly exponential, it resulted in an equal page and link distribution (see Figures 2, 3, 4, 5). Remarkably, the intra-partition ratio (inter-host links inside the same partition) is negligible, as the inter-partition link rate nearly equals to the inter-host ratio. This means that hosts can arbitrarily be shifted from one partition to another one (which is necessary for fast re-balancing with incremental web crawling).

Type	Amount	Percent
Total	601,183,777	100%
Intra-Host	563,992,416	93.81%
Inter-Host	37,191,361	6.19%
Inter-Partition	32,716,628	5.44%
Intra-Partition	4,474,733	0.74%

Table 2. LargeWeb link distribution

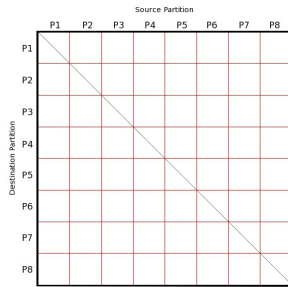


Fig. 2. Partitioned, normalized LargeWeb-Dotplot

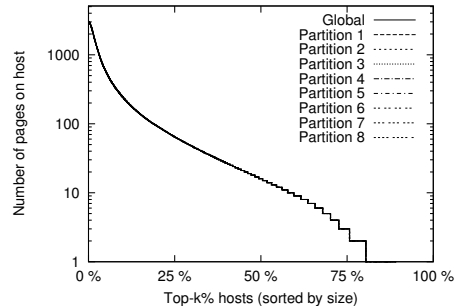


Fig. 3. Partitioned Host Distribution

5.1 Implementation

We have implemented Partitioned PageRank in Java using a P2P-like network with a central coordinator instance. This coordinator is only responsible for arranging the iteration process at partition-level and does not know anything about the rank scores or the link structure. Before the computation, all nodes announce themselves to the coordinator, communicating the hosts they cover. The iteration process is started as soon as all nodes are ready. The coordinator then broadcasts the global host structure to all known nodes and instructs them to iterate. Whenever a node's subgraph changes, it sends lists of external outgoing link targets to the corresponding nodes.

For every iteration step, a node will compute its votes using our reformulated PageRank (Equation 6); the partition itself is again divided in subpartitions processed in parallel. The nodes then aggregate all outgoing inter-partition votes by

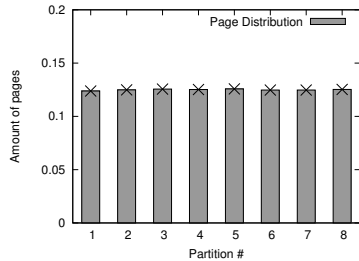


Fig. 4. Pages per Partition

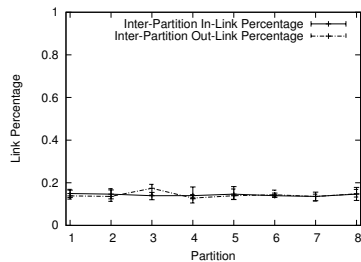


Fig. 5. Inter-partition link source/target distribution

target page and send them directly to the other nodes responsible for these target pages, in the order specified beforehand. Finally, each node reports its local rank status (using the sum and number of its PageRank scores) to the coordinator, in order to compute the global residual δ . As soon as all nodes have succeeded, the coordinator decides whether to continue iterating, by broadcasting another “iterate” command unless the residual reached the defined threshold ϵ .

The addition of new pages during incremental crawling may happen at any time. If the addition covers new hosts, the coordinator selects a node according to the current balancing. From then on, this node is responsible for all pages of that host. The assignment is broadcasted to all nodes in case that there were dangling links to that (previously uncovered) host.

5.2 Results

We conducted most of the experiments on four Linux machines, an AMD Dual Opteron 850 2.4 GHz, 10GB RAM (“A”), an Intel Dual Xeon 2.8 GHz, 6GB RAM (“B”) and two Intel Xeon 3.0 GHz, 1.5GB RAM (“C” and “D”). They were connected via 100MBit Ethernet LAN and not under load before our experiments. We divided the LargeWeb graph into eight partitions and distributed them among the four servers according to available memory (Machine A holds four partitions, B two, C and D one) and performed unbiased PageRank computations.

We examined the convergence behavior, rank distribution and elapsed time both globally and per-partition. All per-partition results matched almost perfectly with the global counterpart and therefore confirmed our assumptions (see Figure 6). The PageRank computation converged below $\epsilon = 10^{-3}$ after 17 iterations, and the entire computation took less than 9 minutes, with only 66 seconds accounted for rank computation, the rest being network I/O. More, with Gigabit-Ethernet and Java Non-Blocking I/O, network communication costs would probably go down to the same magnitude as computation costs.

Compared to the running times of a centralized PageRank computation, our computation is about 10 times faster per iteration, with network I/O included,

and about 75 times faster if network I/O is not counted (when all pages would fit on one massive parallel machine).

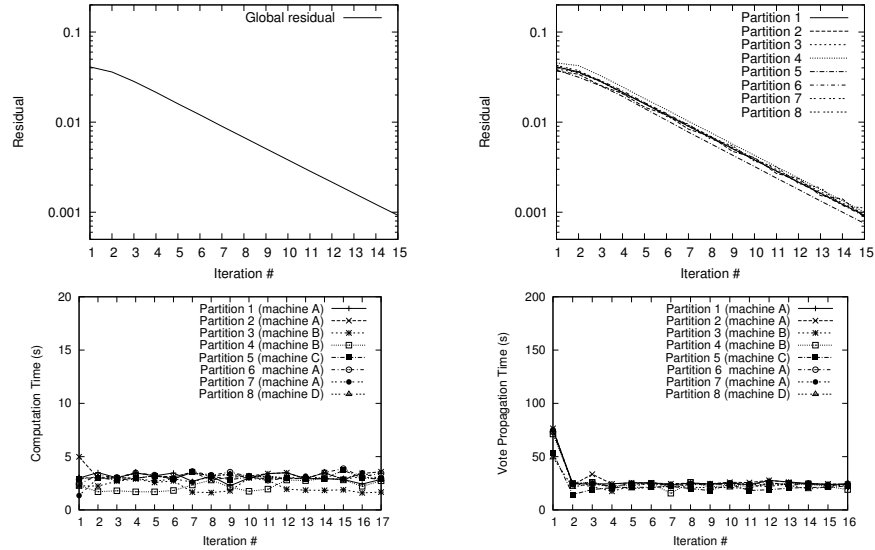


Fig. 6. Partitioned PageRank convergence, vote calculation and network communication times using 8 partitions on 4 machines; $\varepsilon = 0.001$

6 Conclusions and Further Work

We presented an efficient method to perform the PageRank calculation in parallel over arbitrary large web graphs. We accomplished this by introducing a novel two-dimensional view of the web, having the host ID as the only discriminator, as well as by adapting the Gauß-Seidel method for solving linear systems to be used with PageRank. Additionally, we also discussed several other possible applications of our approach, such as a fast re-balancing for incremental crawling (faster than any other approaches known to us). Our next goal is to experiment with several other PageRank specific enhancements, for example with extrapolation methods (that reduce convergence time), under extensive memory demanding scenarios.

References

1. Arvind Arasu, Jasmine Novak, Andrew Tomkins, and John Tomlin. Pagerank computation and the structure of the web: Experiments and algorithms, 2001.
2. Krishna Bharat, Bay-Wei Chang, Monika Rauch Henzinger, and Matthias Ruhl. Who links to whom: Mining linkage between web sites. In *Proc. of the IEEE Intl. Conf. on Data Mining*, pages 51–58, 2001.

3. Sergey Brin, Rajeev Motwani, Lawrence Page, and Terry Winograd. What can you do with a web in your pocket? *Data Engineering Bulletin*, 21(2):37–47, 1998.
4. Andrei Z. Broder, Ronny Lempel, Farzin Maghoul, and Jan Pedersen. Efficient pagerank approximation via graph aggregation. In *Proc. of the 13th International World Wide Web Conference*, pages 484–485, 2004.
5. Yen-Yu Chen, Qingqing Gan, and Torsten Suel. I/o-efficient techniques for computing pagerank, 2002.
6. Junghoo Cho and Hector Garcia-Molina. The evolution of the web and implications for an incremental crawler. In *Proceedings of the Twenty-sixth International Conference on Very Large Databases*, 2000.
7. Nadav Eiron, Kevin S. McCurley, and John A. Tomlin. Ranking the web frontier. In *Proc. of the 13th Intl. Conf. on the World Wide Web*, pages 309–318, 2004.
8. David Gleich, Leonid Zhukov, and Pavel Berkhin. Fast parallel PageRank: A linear system approach. Technical report, Yahoo! Research Labs, 2004.
9. Taher H. Haveliwala. Efficient computation of PageRank. Technical Report 1999-31, Stanford Digital Library Technologies Project, 1999.
10. Sepandar Kamvar, Taher Haveliwala, Christopher Manning, and Gene Golub. Exploiting the block structure of the web for computing PageRank. Technical report, Stanford University, 2003.
11. Sepandar D. Kamvar, Taher H. Haveliwala, Christopher D. Manning, and Gene H. Golub. Extrapolation methods for accelerating PageRank computations. In *Proc. of the 12th Intl. Conf. on the World Wide Web*, pages 261–270, 2003.
12. Sung Jin Kim and Sang Ho Lee. An improved computation of the PageRank algorithm. In *Proc. of the European Conference on Information Retrieval (ECIR)*, pages 73–85, 2002.
13. D. P. Koester, S. Ranka, and G. C. Fox. A parallel gauss-seidel algorithm for sparse power system matrices. In *Proc. of the ACM/IEEE Conf. on Supercomputing*, pages 184–193, 1994.
14. Amy N. Langville and Carl D. Meyer. Deeper inside PageRank, 2004.
15. Chris P. Lee, Gene H. Golub, and Stefanos A. Zenios. A fast two-stage algorithm for computing PageRank. Technical report, Stanford University, 2003.
16. Bundit Manaskasemsak and Arnon Rungsawang. Parallel PageRank computation on a gigabit pc cluster. In *Proceedings of the 18th International Conference on Advanced Information Networking and Application (AINA'04)*, 2004.
17. Netcraft. Web server survey, 2004.
18. Lawrence Page, Sergey Brin, Rajeev Motwani, and Terry Winograd. The PageRank citation ranking: Bringing order to the web. Technical report, Stanford Digital Library Technologies Project, 1998.
19. Karthikeyan Sankaralingam, Simha Sethumadhavan, and James C. Browne. Distributed pagerank for p2p systems. In *Proc. of the 12th IEEE Intl. Symp. on High Performance Distributed Computing (HPDC)*, page 58, 2003.
20. Taher H. Haveliwala Sepandar D. Kamvar and Gene H. Golub. Adaptive methods for the computation of PageRank. Technical report, Stanford University, 2003.
21. Shu-Ming Shi, Jin Yu, Guang-Wen Yang, and Ding-Xing Wang. Distributed page ranking in structured p2p networks. In *Proceedings of the 2003 International Conference on Parallel Processing (ICPP'03)*, pages 179–186, 2003.
22. Taher H. Haveliwala et al. 2001 Crawl of the WebBase project.
23. Yuan Wang and David J. DeWitt. Computing PageRank in a distributed internet search system. In *Proceedings of the 30th VLDB Conference*, 2004.
24. Jie Wu and Karl Aberer. Using SiteRank for P2P Web Retrieval, March 2004.