

UbiCrawler: A Scalable Fully Distributed Web Crawler

Paolo Boldi* Bruno Codenotti† Massimo Santini‡ Sebastiano Vigna§

Abstract

We report our experience in implementing UbiCrawler, a scalable distributed web crawler, using the Java programming language. The main features of UbiCrawler are platform independence, linear scalability, graceful degradation in the presence of faults, a very effective assignment function (based on consistent hashing) for partitioning the domain to crawl, and more in general the complete decentralization of every task. The necessity of handling very large sets of data has highlighted some limitation of the Java APIs, which prompted the authors to partially reimplement them.

1 Introduction

In this paper we present the design and implementation of UbiCrawler, a scalable, fault-tolerant and fully distributed web crawler, and we evaluate its performance both *a priori* and *a posteriori*. The overall structure of the UbiCrawler design was preliminarily described in [2]¹, [5] and [4].

This work is part of a project which aims at gathering large data sets to study the structure of the web. This goes from statistical analysis of specific web domains [3] to estimates of the distribution of classical parameters, such as page rank [20] and to the development of techniques to redesign Arianna, the largest Italian search engine.

Since the first stages of the project, we realized that centralized crawlers are not any longer sufficient to crawl meaningful portions of the web. Indeed, it has been recognized that “as the size of the web grows, it becomes imperative to parallelize the crawling process, in order to finish downloading pages in a reasonable amount of time” [9, 1].

Many commercial and research institutions run their web crawlers to gather data about the web. Even if no code is available, in several cases the basic design has been made public: this is the case, for instance, of Mercator [18] (the Altavista crawler), of the original Google crawler [6], and of some crawlers developed within the academic community [23, 24, 22].

Nonetheless, little published work actually investigates the fundamental issues underlying the parallelization of the different tasks involved in the crawling process. In particular, all approaches we are aware of employ some kind of centralized manager that decides which URLs are to be visited, and that stores the URLs which have already been crawled. At best, these components can be replicated and their work can be partitioned statically.

In contrast, when designing UbiCrawler, we have decided to decentralize every task, with obvious advantages in terms of scalability and fault tolerance.

Essential features of UbiCrawler are

*Dipartimento di Scienze dell’Informazione, Università degli Studi di Milano, via Comelico 39/41, I-20135 Milano, Italy. boldi@dsi.unimi.it

†Department of Computer Science, The University of Iowa, 14 Maclean Hall, Iowa City IA 52240 (on leave from IIT-CNR, Pisa, Italy) bcodenot@cs.uiowa.edu

‡Dipartimento di Scienze Sociali, Cognitive e Quantitative, Università di Modena e Reggio Emilia, via Fratelli Manfredi I-42100 Reggio Emilia, Italy. msantini@unimo.it

§Contact author. Dipartimento di Scienze dell’Informazione, Università degli Studi di Milano, via Comelico 39/41, I-20135 Milano, Italy. vigna@acm.org. Phone: +39-0258356324. Fax: +39-0258356373.

¹At the time, the name of the crawler was *Trovatore*, later changed into UbiCrawler when the authors learned about the existence of an Italian search engine named Trovatore.

- platform independence;
- full distribution of every task (no single point of failure and no centralized coordination at all);
- locally computable URL assignment based on consistent hashing;
- tolerance to failures: permanent as well as transient failures are dealt with gracefully;
- scalability.

As we will outline in Section 2, these characteristics are the byproduct of a well defined design goal: fault tolerance and full distribution (lack of any centralized control) are assumptions which have guided our architectural choices. For instance, while there are several reasonable ways to partition the domain to be crawled if we assume the presence of a central server, it becomes harder to find an assignment of URLs to different agents which is fully distributed, does not require too much coordination, and allows us to cope with failures.

In Section 2 we present the overall design of UbiCrawler discussing in particular the requirements which guided our choices. Section 3 gives a high level description of the software architecture of the crawler and Section 4 introduces the assignment function used to distribute the URLs to be crawled, and gives some general results about its properties. The implementation issues faced in developing UbiCrawler are detailed in Section 5, while Section 6 is devoted to the performance evaluation, both from an analytical and an empirical point of view. Finally, Section 7 contrasts our results with related work in the literature.

2 Design Assumptions, Requirements, and Goals

In this section we give a brief presentation of the most important design choices which have guided the implementation of UbiCrawler. More precisely, we sketch general design goals and requirements, as well as assumptions on the type of faults that should be tolerated.

Full distribution. In order to achieve significant advantages in terms of programming, deployment, and debugging, a parallel and distributed crawler should be composed of identically programmed *agents*, distinguished by a unique identifier only. This has a fundamental consequence: each task must be performed in a fully distributed fashion, that is, no central coordinator can exist. Full distribution is instrumental in obtaining a scalable, easily configurable system that has no single point of failure.

We also do not want to rely on any assumption concerning the location of the agents, and this implies that latency can become an issue, so that we should minimize communication to reduce it.

Balanced locally computable assignment. The distribution of URLs to agents is an important problem, crucially related to the efficiency of the distributed crawling process.

We identify the three following goals:

- At any time, each URL should be assigned to a specific agent, which is the only *responsible* for it, to avoid undesired data replication.
- For any given URL, the knowledge of its responsible agent should be locally available. In other words, every agent should have the capability to compute the identifier of the agent responsible for a URL, without communication. This feature reduces the amount of inter-agent communication; moreover, if an agent detects a fault while trying to assign a URL to another agent, it will be able to choose the new responsible without further communication.
- The distribution of URLs should be *balanced*, that is, each agent should be responsible for approximately the same number of URLs. In case of heterogeneous agents, the number of URLs should be proportional to the agent's available resources (such as memory, hard disk capacity etc.).

Scalability. The number of pages crawled per second and agent should be (almost) independent of the number of agents. In other words, we expect the throughput to grow linearly with the number of agents.

Politeness. A parallel crawler should never try to fetch more than one page at a time from a given host. Moreover, a suitable delay should be introduced between two subsequent requests to the same host.

Fault tolerance. A distributed crawler should continue to work under *crash faults*, that is, when some agents abruptly die. No behavior can be assumed in the presence of this kind of crash, except that the faulty agent stops communicating; in particular, one cannot prescribe any action to a crashing agent, or recover its state afterward². When an agent crashes, the remaining agents should continue to satisfy the “Balanced locally computable assignment” requirement: this means, in particular, that URLs of the crashed agent will have to be redistributed.

This has two important consequences:

- It is not possible to assume that URLs are statically distributed.
- Since the “Balanced locally computable assignment” requirement must be satisfied *at any time*, it is not reasonable to rely on a distributed reassignment protocol after a crash. Indeed, during the reassignment the requirement would be violated.

3 The Software Architecture

UbiCrawler is composed by several agents that autonomously coordinate their behavior in such a way that each of them scans its share of the web. An agent performs its task by running several threads, each dedicated to the visit of a single host. More precisely, each thread scans a single host using a breadth-first visit. We make sure that different threads visit different hosts at the same time, so that each host is not overloaded by too many requests. The outlinks that are not local to the given host are dispatched to the right agent, which puts them in the queue of pages to be visited. Thus, the overall visit of the web is breadth first, but as soon as a new host is met, it is entirely visited (possibly with bounds on the depth reached or on the overall number of pages), again in a breadth-first fashion.

More sophisticated approaches (which can take into account suitable priorities related to URLs, such as, for instance, their rank) can be easily implemented. However it is worth noting that several authors (see, e.g., [19]) have argued that breadth-first visits tends to find high quality pages early on in the crawl. A deeper discussion about page quality is given in Section 6.

An important advantage of per-host breadth-first visits is that DNS requests are infrequent. Web crawlers that use a global breadth-first strategy must work around the high latency of DNS servers: this is usually obtained by buffering requests through a multithreaded cache. Similarly, no caching is needed for the `robots.txt` file required by the “Robot Exclusion Standard” [16]; indeed such file can be downloaded when a host visit begins.

Assignment of hosts to agents takes into account the mass storage resources and bandwidth available at each agent. This is currently done by means of a single indicator, called *capacity*, which acts as a weight used by the assignment function to distribute hosts. Under certain circumstances, each agent a gets a fraction of hosts proportional to its capacity C_a (see Section 4 for a precise description of how this works). Note that even if the number of URLs per host varies wildly, the distribution of URLs among agents tends to even out during large crawls. Besides empirical statistical reasons for this, there are also other motivations, such as the usage of policies for bounding the maximum number of pages crawled from a host and the maximum depth of a visit. Such policies are necessary to avoid (possibly malicious) *web traps*.

²Note that this is radically different from milder assumptions, as for instance saying that the state of a faulty agent can be recovered. In the latter case, one can try to “mend” the crawler’s global state by analyzing the state of the crashed agent.

Finally, an essential component of UbiCrawler is a *reliable failure detector* [8], that uses timeouts to detect crashed agents; reliability refers to the fact that a crashed agent will eventually be distrusted by every active agent (a property that is usually referred to as *strong completeness* in the theory of failure detectors). The failure detector is the only synchronous component of UbiCrawler (i.e., the only component using timings for its functioning); all other components interact in a completely asynchronous way.

4 The Assignment Function

In this section we describe the assignment function used by UbiCrawler, and we explain why this function makes it possible to decentralize every task and to achieve our fault-tolerance goals.

Let \mathcal{A} be our set of agent identifiers (i.e., potential agent names), and $\mathcal{L} \subseteq \mathcal{A}$ be the set of alive agents: we have to assign hosts to agents in \mathcal{L} . More precisely, we have to set up a function δ that, for each nonempty set \mathcal{L} of alive agents, and for each host h , delegates the responsibility of fetching (URLs from) h to the agent $\delta_{\mathcal{L}}(h) \in \mathcal{L}$.

The following properties are desirable for an assignment function:

1. *Balancing*. Each agent should get approximately the same number of hosts; in other words, if m is the (total) number of hosts, we want that $|\delta_{\mathcal{L}}^{-1}(a)| \sim m/|\mathcal{L}|$ for each $a \in \mathcal{L}$.
2. *Contravariance*. The set of hosts assigned to an agent should change in a contravariant manner with respect to the set of alive agents across a deactivation and reactivation. More precisely, if $\mathcal{L} \subseteq \mathcal{L}'$ then $\delta_{\mathcal{L}}^{-1}(a) \supseteq \delta_{\mathcal{L}'}^{-1}(a)$; that is to say, if the number of agents grows, the portion of the web crawled by each agent must shrink. Contravariance has a fundamental consequence: if a new set of agents is added, no old agent will ever lose an assignment in favor of another old agent; more precisely, if $\mathcal{L} \subseteq \mathcal{L}'$ and $\delta_{\mathcal{L}'}(h) \in \mathcal{L}$ then $\delta_{\mathcal{L}'}(h) = \delta_{\mathcal{L}}(h)$; this guarantees that at any time the set of agents can be enlarged with minimal interference with the current host assignment.

Note that satisfying partially the above requirement is not difficult: for instance, a typical approach used in non-fault-tolerant distributed crawlers is to compute a modulo-based hash function of the host name. This has very good balancing properties (each agent gets approximately the same number of hosts), and certainly can be computed locally by each agent knowing just the set of alive agents.

However, what happens when an agent crashes? The assignment function can be computed again, giving however a different result for almost all hosts. The *size* of the sets of hosts assigned to each agent would grow or shrink contravariantly, but the *content* of those sets would change in a completely chaotic way. As a consequence, after a crash most pages will be stored by an agent that should not have fetched them, and they could mistakenly be re-fetched several times³.

Clearly, if a central coordinator is available or if the agents can engage a kind of “resynchronization phase” they could gather other information and use other mechanisms to redistribute the hosts to crawl. However, we would have just shifted the fault-tolerance problem to the resynchronization phase—faults in the latter would be fatal.

4.1 Background

Although it is not completely obvious, it is not difficult to show that contravariance implies that each possible host induces a total order (i.e., a permutation) on \mathcal{A} ; more precisely, a contravariant assignment is equivalent to a function that assigns an element of $S_{\mathcal{A}}$ (the symmetric group over \mathcal{A} , i.e., the set of all permutations elements of \mathcal{A} , or equivalently, the set of all total orderings of elements

³For the same reason, a modulo-based hash function would make it difficult to increase the number of agents during a crawl.

of \mathcal{A}) to each host: then, $\delta_{\mathcal{L}}(h)$ is computed by taking, in the permutation associated to h , the first agent that belongs to the set \mathcal{L} .

A simple technique to obtain a balanced, contravariant assignment function consists in trying to generate such permutations, for instance, using some bits extracted from a host name to seed a (pseudo)random generator, and then permuting randomly the set of possible agents. This solution has the big disadvantage of running in time and space proportional to the set of possible agents (which one wants to keep as large as feasible). Thus, we need a more sophisticated approach.

4.2 Consistent Hashing

Recently, a new hashing technique called *consistent hashing* [14, 15] has been proposed for the implementation of a system of distributed web caches (a different approach to the same problem can be found in [10]). The idea of consistent hashing is very simple, yet profound.

As we noted, for a typical hash function, adding a bucket (i.e., a new place in the hash table) is a catastrophic event. In consistent hashing, instead, each bucket is replicated a fixed number κ of times, and each copy (we shall call it a *replica*) is mapped randomly on the unit circle. When we want to hash a key, we compute in some way from the key a point in the unit circle, and find its nearest replica: the corresponding bucket is our hash. The reader is referred to [14] for a detailed report on the powerful features of consistent hashing, which in particular give us balancing for free. Contravariance is also easily verified.

In our case, buckets are agents, and keys are hosts. We must be very careful, however, if we want the contravariance (2) to hold, because mapping randomly the replicas to the unit circle each time an agent is started will not work; indeed, δ would depend not only on \mathcal{L} , but also on the choice of the replicas. Thus, *all* agents should compute the same set of replicas corresponding to a given agent, so that, once a host is turned into a point of the unit circle, all agents will agree on who is responsible for that host.

4.3 Identifier–Seeded Consistent Hashing

A method to fix the set of replicas associated to an agent and try to maintain the good randomness properties of consistent hashing is to derive the set of replicas from a very good random number generator seeded with the agent identifier: we call this approach *identifier-seeded consistent hashing*. We have opted for the Mersenne Twister [17], a fast random generator with an extremely long cycle that passes very strong statistical tests.

However this solution imposes further constraints: since replicas cannot overlap, any discretization of the unit circle will incur in the Birthday paradox—even with a very large number of points, the probability that two replicas overlap will become non-negligible. Indeed, when a new agent is started, its identifier is used to generate the replicas for the agent. However, if during this process we generate a replica that is already assigned to some other agent, we must force the new agent to choose another identifier.

This solution might be a source of problems if an agent goes down for a while and discovers a conflict when it is restarted. Nonetheless, some standard probability arguments show that with a 64-bit representation for the elements of the unit circle there is room for 10^4 agents with a conflict probability of 10^{-12} .

We remark that a theoretical analysis of the balancing produced by identifier-seeded consistent hashing is most difficult, if not impossible (unless, of course, one uses the working assumption that replicas behave as if they were randomly distributed). Thus, we report experimental data: in Figure 1 we can see that once a substantial number of hosts have been crawled, the deviation from perfect balancing is less than 6% for small as well as for large sets of agents when $\kappa = 100$, that is, we use 100 replicas per bucket (thin lines); if $\kappa = 200$, the deviation decreases to 4.5% (thick lines).

We have implemented consistent hashing as follows: the unit interval can be mapped on the whole set of representable integers, and then replicas can be kept in a balanced tree whose keys are integers. This allows us to hash a host in logarithmic time (in the number of alive agents). By keeping

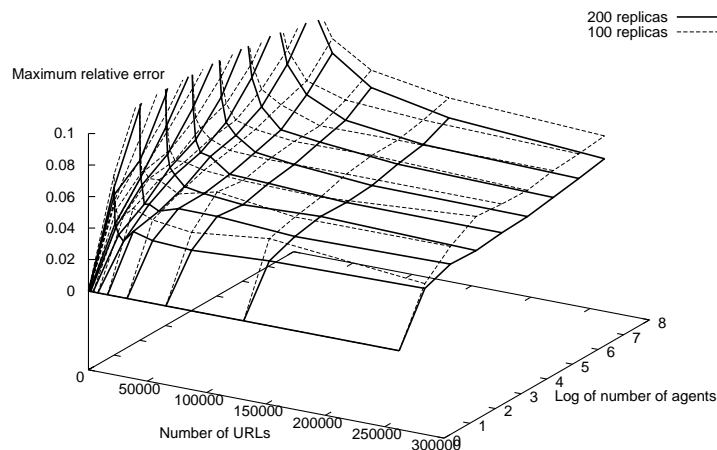


Figure 1: Experimental data on identifier-seeded consistent hashing. Deviation from perfect balancing is less than 6% with 100 replicas (thin lines), and less than 4.5% with 200 replicas (thick lines).

the leaves of the tree in a doubly linked chain we can also easily implement the search for the next nearest replica.

As we already mentioned, an important feature of UbiCrawler is that it can run on heterogeneous hardware, with different amount of available resources. To this purpose, an agent is associated with a number of replicas proportional to its capacity, and this guarantees that the assignment function distributes hosts evenly with respect to the mass storage available at each agent.

Moreover, the number of threads of execution for each agent can be tuned to suit network bandwidth or CPU limitations. Note however that an excessive number of threads can lead to contention on shared data structures, such as the agent's store, and to excessive CPU load, with corresponding performance degradation.

5 Implementation issues

As we already mentioned, several key ideas in web crawling have been made public and discussed in many seminal papers. UbiCrawler builds on this knowledge and uses ideas from previous crawlers, such as Rabin's fingerprinting technique [18].

We decided to develop UbiCrawler as a pure 100% Java application. The choice of Java™ 2 as implementation language is mainly motivated by our need to achieve platform-independence, a necessity that is especially urgent for a fully distributed P2P-like application. Currently UbiCrawler consists of about one-hundred twenty Java classes and interfaces organized in fifteen packages, with about 800 methods and more than 12 000 lines of code.

Of course, Java imposes a certain system overhead, when contrasted with a C/C++ implementation. Nevertheless, our tests show that the speed of UbiCrawler is limited by network bandwidth, and not by CPU power. In fact, the performance penalty of Java is much smaller than usually believed; for instance, the implementors of the CERN Colt package [13] claim a 2.5 linear performance penalty against *hand-crafted assembler*. The (realistic) user perception of an intrinsic Java slowness is mainly due to the bad performance of the Swing window toolkit.

Moreover, Java made it possible to adopt *Remote Method Invocation* [21], a technology that enables one to create distributed applications in which the methods of remote Java objects can be invoked from other Java virtual machines (possibly on different hosts), using object serialization to implicitly marshal and unmarshal parameters. This freed us from the necessity of implementing

communication protocols among agents.

The components of each agent interact as semi-independent modules, each running possibly more than one thread. To bound the amount of information exchanged on the network, each agent is confined to live in a single machine. Nonetheless, different agents may (and typically do) run on different machines, and interact using RMI.

The intensive use of Java APIs in a highly distributed and time/space critical project has highlighted some limitations and issues that led us to devise new ad-hoc solutions, some of which turned out to be interesting *per se*.

Space/time-efficient type-specific collections. The `Collection` and `Map` hierarchies in the `java.util` package are a basic tool that is most useful in code development. Unfortunately, because of the awkward management of primitive types (to be stored in collection they need to be wrapped in suitable objects) those hierarchies are not suitable for handling primitive types, a situation that often happens in practice. If you need a set of integers, you should wrap every single integer into an `Integer` object. Apart for space inefficiency, object creation is a highly time-consuming task, and the creation of many small objects makes garbage collection problematic, a fact that becomes dramatic in a distributed setting, where responsiveness is critical.

More issues derive from the way collections and maps are implemented in the standard API. For example, a `HashMap` is realized using closed addressing, so every entry in the table has an additional reference, and moreover it caches hash codes; hence, an entry with a pair of `int` (that should minimally take 8 bytes) requires the allocation of 3 objects (two `Integer` objects for wrapping the two integers, and an entry object), and the entry contains three references (key, value and next field) and an additional integer field to keep the hash code cached. A `HashSet` is implemented as a single-valued `HashMap`.

Each `UbiCrawler` agent keeps track of the URLs it has visited: this is obtained via a hash table that stores 64-bit CRCs (a.k.a. fingerprints) of the URLs. Indeed, this table turns out to be responsible for most of the memory occupancy. Storing this table using the standard APIs would reduce by a factor of at least 20 the number crawlable URLs (even worse, the number of objects in memory would make garbage collections so time consuming to produce timeouts in inter-process communications).

All these considerations led to the development of a package providing alternatives to the sets and maps defined in `java.util`. This package, named `fastUtil`, contains 537 classes that offer type-specific mappings (such as `Int2LongOpenHashMap`). They all implement the standard Java interfaces, but offer also polymorphic methods for easier access and reduced object creation. The algorithmic techniques used in our implementation are rather different than those of the standard API (e.g., open addressing, threaded balanced trees with bidirectional iterators, etc.), and provide the kind of performance and controlled object creations that we needed. These classes have been released under the GNU Lesser General Public License.

Robust, fast, error-tolerant HTML parsing. Every crawling thread, after fetching a page, needs to parse it before storing; parsing is required both to extract hyperlinks that are necessary for the crawling to proceed, and to obtain other relevant information (e.g., the charset used in the page, in case it differs from the one specified in its headers; the set of words contained in the page, an information that is needed for indexing, unless one wants to make this analysis off-line with a further parsing step). The current version of `UbiCrawler` uses a highly optimized HTML/XHTML parser that is able to work around most common errors. On a standard PC, performance is about 600 page/s (this includes URL parsing and word occurrence extraction).

String and StringBuffer. The Java string classes are a well-known cause of inefficiency. In particular, `StringBuffer` is synchronized, which implies a huge performance hit in a multithreaded application. Even worse, `StringBuffer` has equality defined by reference (i.e., two buffers with the same content are not equal), so even a trivial task such as extracting word occurrences and storing them in a data structure poses nontrivial problems. In the end, we rewrote a string

class lying halfway between `String` and `StringBuffer`. The same problems have been reported by the authors of Mercator [12], who also claim to have rewritten the Java string classes.

6 Performance Evaluation

The goal of this section is to discuss UbiCrawler in the framework of the classification given in [9], and to analyze its scalability and fault-tolerance features. In particular, we consider the most important properties identified by [9] (degree of distribution, coordination, partitioning techniques, coverage, overlap, and communication overhead) and contrast UbiCrawler against them.

Degree of distribution. A parallel crawler can be intra-site, or distributed, that is, its agents can communicate either through a LAN or through a WAN. UbiCrawler is a distributed crawler which can run on any kind of network.

Coordination. In the classification of [9], agents can use a different amount of coordination: at one extreme, all agents crawl the network independently, and one hopes that their overlap will be small due to a careful choice of the starting URLs; at the other extreme, a central coordinator divides the network either *statically* (i.e., before the agents actually start) or *dynamically* (i.e., during the crawl).

As for UbiCrawler, the assignment function gives rise to a kind of coordination that does not fit the models and the options suggested above. Indeed, the coordination is dynamic, but there is no central authority that handles it. Thus, in a sense, all agents run independently, but they are at the same time tightly and distributedly coordinated. We call this feature *distributed dynamic coordination*.

Partitioning techniques. The web can be partitioned in several ways; in particular, the partition can be obtained from URL-based hash, host-based hash or hierarchically, using, for instance, Internet domains. Currently, UbiCrawler uses a host-based hash; note that since [9] does not consider consistent hashing, some of the arguments about the shortcomings of hashing functions are no longer true for UbiCrawler.

Coverage. It is defined as $\frac{c}{u}$, where c is the number of actually crawled pages, and u the number of pages the crawler as a whole had to visit.

If no faults occur, UbiCrawler achieves coverage 1, which is optimal. Otherwise, it is in principle possible that some URLs that were stored locally by a crashed agent will not be crawled. However, if these URLs are reached along other paths after the crash, they will clearly be fetched by the new agent responsible for them.

Overlap. It is defined as $\frac{n-u}{u}$, where n is the total number of pages crawled by *alive* agents and u the number of *unique* pages; note that $u < n$ can happen if the same page has been erroneously fetched several times.

Even in the presence of crash faults, UbiCrawler achieves overlap 0, which is optimal. However, if we consider transient faults, where an agent may be temporarily unavailable, we cannot guarantee the absence of duplications. In particular, we cannot prevent other agents from fetching a URL that a temporarily unavailable agent already stores, because we cannot foresee whether the fault is transient or not (unless, of course, we accept a potentially incomplete coverage).

Nevertheless, note that after a transient fault UbiCrawler autonomously tries to converge to a state with overlap 0 (see Section 6.1.1). This property is usually known as *self-stabilization*, a technique for protocol design introduced by Dijkstra [11].

Communication overhead. It is defined as $\frac{e}{n}$, where e is the number of URLs exchanged by the agents during the crawl and n is the number of crawled pages.

Assuming that every page contains λ links to other sites (on average), we have that n crawled pages will give rise to λn URLs that must be potentially communicated to other agents⁴. Due to the balancing property of the assignment function, at most

$$\lambda n \frac{\sum_{a \neq \bar{a}} C_a}{\sum_a C_a} < \lambda n$$

messages will be sent across the network, where a ranges in the set of alive agents, and \bar{a} is the agent that fetched the page (recall that C_a is the capacity of agent a). By the definition of [9], our communication overhead is thus less than λ . It is an interesting feature that the number of messages is *independent of the number of agents*, and depends only on the number of crawled pages and on λ . In other words, a large number of agents will generate more network traffic, but this is due to the fact that they are fetching more pages, and not to a design bottleneck.

Quality. It is a complex measure of ‘importance’ or ‘relevance’ of crawled pages as determined by suitable ranking techniques; an important challenge is to build a crawler that tends to collect high-quality pages during the early stages of the crawling process.

As we already mentioned, currently UbiCrawler uses a parallel per-host breadth-first visit, without dealing with ranking and quality-of-page issues. This is because our immediate goal is to focus on scalability of the crawler itself and on the analysis of some portions of the web, as opposed to building a search engine. Nonetheless, since a breadth-first single-process visit tends to visit high-quality pages first [19], it is natural to ask whether our strategy works well or not.⁵

To be more precise, UbiCrawler has a limit on the depth of any host visit. Once the limit is reached, the visit terminates. In particular, this means that by setting the limit to 0, UbiCrawler performs a pure breadth-first visit, whereas by setting the limit to higher values, the visit resembles more and more a depth-first one.

Figure 2 shows the cumulative PageRank during a crawl of about 45,000,000 pages of the domain .it; We compare (from top to bottom) an ideal omniscient strategy that visits pages of high PageRank first; then, a breadth-first visit, the actual UbiCrawler visit and a depth-first visit. As the crawl was performed with a high depth limit (8), the UbiCrawler visit is nearer to the results of the depth-first visit.

6.1 Fault Tolerance

To the best of our knowledge, no commonly accepted metrics exist for estimating the fault tolerance of distributed crawlers, since the issue of faults has not been taken into serious consideration up to now. It is indeed an interesting and open problem to define a set of measures to test the robustness of parallel crawlers in the presence of faults. Thus, we give an overview of the reaction of UbiCrawler agents to faults.

UbiCrawler agents can die or become unreachable either expectedly (for instance, for maintenance) or unexpectedly (for instance, because of a network problem). At any time, each agent has its own view of which agents are alive and reachable, and these views do not necessarily coincide.

Whenever an agent dies abruptly, the failure detector discovers that something bad has happened (e.g., using timeouts). Thanks to the properties of the assignment function, the fact that different agents have different views of the set of alive agents does not disturb the crawling process. Suppose, for instance, that a knows that b is dead, whereas a' does not. Because of contravariance, the only

⁴Note that in principle not all URLs must be necessarily communicated to other agents; one could just rely on the choice of a good seed to guarantee that no pages will be lost. Nonetheless, in a worst-case scenario, to obtain coverage 1 all URLs not crawled *must* be communicated to some other agent.

⁵Of course, it will be possible to order pages according to a ranking function, using, for instance, backlink information, at a later stage of this project.

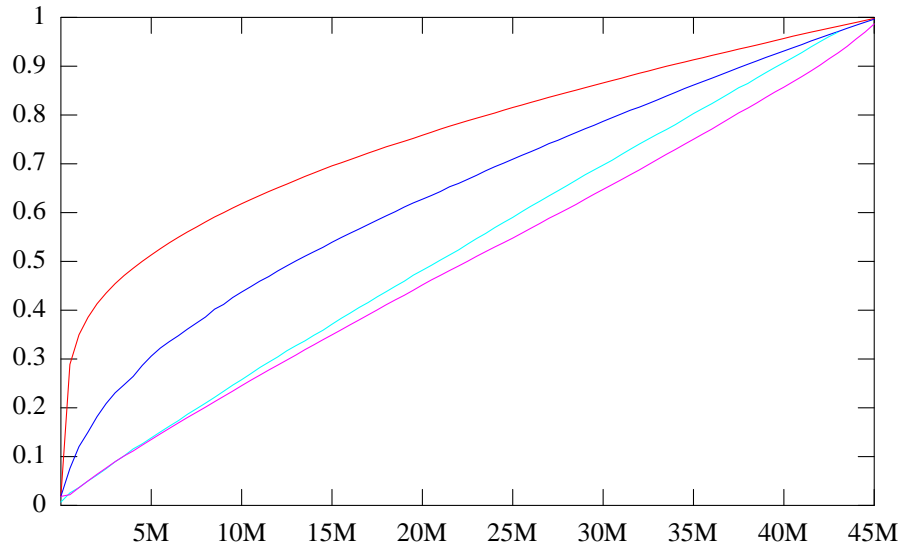


Figure 2: Cumulative PageRank of crawled pages as a function of the number of crawled URLs.

difference between a and a' in assignments of host to agents is the set of hosts pertaining to b . Agent a correctly dispatches these hosts to other agents, and agent a' will do the same as soon as it realizes that b is dead, which will happen, in the worst case, when it tries to dispatch a URL to b . At this point, b will be believed dead, and the host dispatched correctly. Thus, a and a' will never dispatch hosts to different agents.

Another consequence of this design choice is that agents can be dynamically added during a crawl, and after a while all pages for which they are responsible will be removed from the stores of the agents that fetched them before the new agent's birth. In other words, making UbiCrawler self-stabilizing by design gives us not only fault tolerance, but also a greater adaptivity to dynamical configuration changes.

6.1.1 Page Recovery

An interesting feature of contravariant assignment functions is that they allow to guess easily who could have fetched previously a page for which an agent is responsible in the present configuration. Indeed, if a is responsible for the host h , then the agent responsible for h before a was started is the one associated to the next-nearest replica. This allows us to implement a *page recovery protocol* in a very simple way. Under certain conditions, the protocol allows to avoid re-fetching several times the same page even in the presence of faults.

The system is parametrized by an integer t : each time an agent is going to fetch a page of a host for which it is currently responsible, it first checks whether the next-nearest t agents have already fetched that page. It is not difficult to prove that this guarantees page recovery as long as *no more than t agents were started since the page was crawled*. Note that the number of agents that crashed is completely irrelevant.

This approach implies that if we want to accept t (possibly transient) faults without generating overlap (except for the unavoidable cases discussed in Section 6), we have to increase by a linear factor of t the network traffic, as any fetched page will generate at least t communications. This is not unreasonable: typically, in a distributed system the number of rounds required to solve a problem (for instance, consensus) is linearly related to the maximum number of faults.

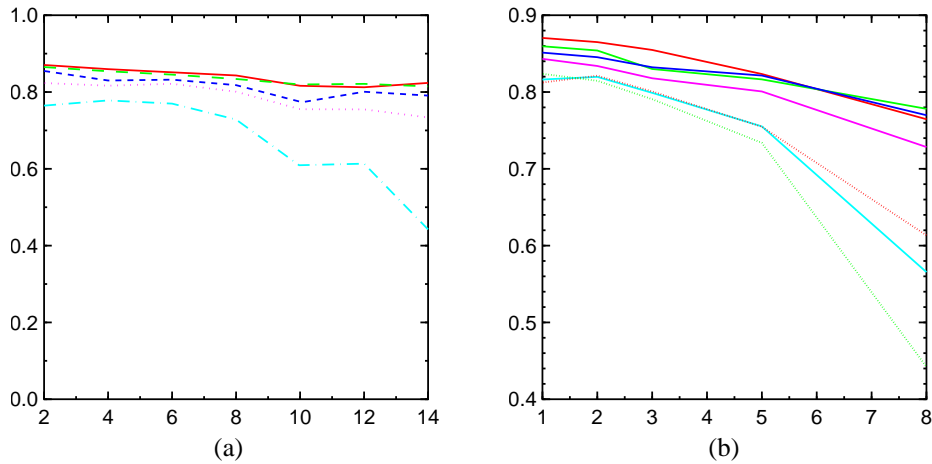


Figure 3: Average number of pages crawled per second and thread, that is, work per thread. Graph (a) shows how work changes when the number of agents changes; the different patterns represent the number of threads (solid line=1, long dashes=2, short dashes=3, dots=5, dash-dots=8). Graph (b) shows how work changes when the number of threads changes; the different lines represent a different number of agents (from 2 to 14, higher to lower).

6.2 Scalability

In a highly scalable system, one should guarantee that the work performed by every thread is constant as the number of threads changes, i.e., that the system and communication overheads do not reduce the performance of each thread. The figures given in Section 6 show that the amount of network communication grows linearly with the number of downloaded pages, a fact which implies that the performance of each UbiCrawler thread is essentially independent of the number of agents. We have measured how the average number of pages stored per second and thread changes when the number of agents, or the number of threads per agent, changes. For example, Figure 3 plots the resulting data for the African domain (these data were gathered during the crawl used for [3]).

Graph (a) shows how work per thread changes when the number of agents increases. In a perfectly scalable system, all lines should be horizontal (which would mean that by increasing the number of agents we could arbitrarily accelerate the crawling process). There is a slight drop in the second part of the first graph, that becomes significant with eight threads. The drop in work, however, is in this case an artifact of the test, caused by our current limitations in terms of hardware resources: to run experiments using more than seven agents, we had to start two agents per machine, and the existence of so many active processes unbearably raised the CPU load, and led to hard disk thrashing. We have decided to include anyway the data because they show almost constant work for a smaller number of threads and for less than eight agents.

Graph (b) shows how work per thread changes when the number of threads per agent increases. In this case, data contention, CPU load and disk thrashing become serious issues, and thus the work performed by each single thread reduces. The drop in work, however, is strongly dependent on the hardware architecture and, again, the reader should take with a grain of salt the lower lines, which manifest the artifact already seen in graph (a).

Just to make these graphs into actual figures, note that a system with sixteen agents, each running four threads, can fetch about 4 500 000 pages a day, and we expect these figures to scale almost linearly with the number of agents, if sufficient network bandwidth is available.

The tests above have been run using a network simulation that essentially provided infinite bandwidth. When network latency is involved, the number of threads can be raised to much higher figures, as thread contention is greatly reduced (and replaced by waiting for the network to provide data). In real crawls, using 50 or more threads, UbiCrawler can download more than 10 000 000 pages per day using five 1GHz PCs (at this point our link was saturated, so we could not try to increase parallelism).

Again, this includes the precomputation of the list of word occurrences of each page.

7 Related works

Although, as mentioned in the introduction, most details about the design and implementation issues of commercial crawlers are not public, there are some highly performant, scalable crawling systems whose structure has been described and discussed by the authors; among them, two distributed crawlers that might be compared to UbiCrawler are *Mercator* [18], used by AltaVista, the spider discussed in [22], and the crawler presented in [23].

Mercator is a high-performance web crawler whose components are loosely coupled; indeed, they can be distributed across several computing units. However, there is a central element, the *frontier*, which keeps track of all the URLs that have been crawled up to now and that filters new requests.

In the original description of *Mercator* this component was unique and centralized. Recently, the authors have added the possibility of structuring a crawler as a *hive*: hosts are statically partitioned among a finite number of drones (with independent crawling and analysis components). However, this does not address the main problem, that is, that all the information about the set of URLs that have been crawled is centralized in the frontier component of each drone. Indeed, *Mercator* uses a very ingenious mix of Rabin fingerprinting and compressed hash tables to access these sets efficiently. On the contrary, UbiCrawler spreads dynamically and evenly among all agents this information.

On the other hand, *Mercator* has a much more complete content handling, providing several protocol modules (Gopher, ftp, etc.) and, more importantly, a *content-seen* module that filters URLs with the same content as URLs that have already been crawled (it should be noted, however, that the authors do not explain how to implement a cross-drone content-seen module).

The spider discussed in [22] is developed using C++ and Python, and the various components interact using socket connections for small message exchanges, and NFS (Network File System) for large messages.

The downloading components communicate with two central components, called *crawl manager* and *crawl application*. The crawl application is responsible for parsing downloaded pages, compressing and storing them; the application is also in charge of deciding the visit policy. The crawl application communicates the URL to be crawled to the manager, that then dispatches the URLs to the downloaders; the manager takes care of issues such as robot-exclusion, speed-rate control, DNS resolution etc.

The described architecture cannot be scaled to an arbitrary number of downloaders, though: the presence of a centralized parser and dispatcher are a bottleneck. The authors solve this problem by partitioning the set of URLs *statically* into k classes, and then using k crawl applications, each responsible for the URLs in one of the classes: the technique adopted here is similar to that of the Internet Archive crawler [7]. The downloaders, thus, communicate each page to the application responsible for that URL. The number of crawl manager used can be reduced by connecting more applications to the same manager. It is worth mentioning that, since the assignment of URLs to applications is fixed statically, the number and structure of crawl applications cannot be changed during runtime (even though one may change the set of downloaders).

The set of visited URLs, maintained by each crawl application, is kept partly in the main memory (using a balanced tree) and partly on disk. Polite crawling is implemented using a domain-based throttling technique that scrambles the URLs in random order; of course, we do not need such technique, because no thread is allowed to issue requests to a host that is currently being visited by another thread.

A notable exception to the previous cases is described in [23], where the authors propose solutions for a completely dynamic distribution of URLs by means of a two-stage URL-distribution process: first of all URLs are mapped to a large array containing agent identifiers; then, the agent obtained from the array has responsibility for the URL. The entries of the array, indeed, act much like replicas in consistent hashing.

However, there are two major drawbacks: first of all, the authors do not explain how to manage births or deaths of *more than one agent*. The technique of array renumbering given in the paper is *not* guaranteed to give a balanced assignment after a few renumbering; moreover, there is no guarantee that if the same agent dies for a short time and then gets alive again it will get the same URL assignment (i.e., contravariance), which is one of the main features of consistent hashing.

8 Conclusions

We have presented UbiCrawler, a fully distributed, scalable and fault-tolerant web crawler. We believe that UbiCrawler introduces new ideas in parallel crawling, in particular the use of consistent hashing as a mean to completely decentralize the coordination logic, graceful degradation in the presence of faults, and linear scalability.

The development of UbiCrawler highlighted also some weaknesses of the Java API, which we have been able to overcome by using, when necessary, better algorithms.

UbiCrawler is an ongoing project, and our current goal is to test the crawler on larger and larger portions of the Web.

References

- [1] Arvind Arasu, Junghoo Cho, Hector Garcia-Molina, Andreas Paepcke, and Sriram Raghavan. Searching the web. *ACM Transactions on Internet Technology*, 1(1):2–43, 2001.
- [2] Paolo Boldi, Bruno Codenotti, Massimo Santini, and Sebastiano Vigna. Trovatore: Towards a highly scalable distributed web crawler. In *Poster Proc. of Tenth International World Wide Web Conference*, pages 140–141, Hong Kong, China, 2001. Winner of the Best Poster Award.
- [3] Paolo Boldi, Bruno Codenotti, Massimo Santini, and Sebastiano Vigna. Structural properties of the African web. In *Poster Proc. of Eleventh International World Wide Web Conference*, Honolulu, USA, 2002.
- [4] Paolo Boldi, Bruno Codenotti, Massimo Santini, and Sebastiano Vigna. UbiCrawler: A scalable fully distributed web crawler. In *Proc. AusWeb02. The Eighth Australian World Wide Web Conference*, 2002.
- [5] Paolo Boldi, Bruno Codenotti, Massimo Santini, and Sebastiano Vigna. UbiCrawler: Scalability and fault-tolerance issues. In *Poster Proc. of Eleventh International World Wide Web Conference*, Honolulu, USA, 2002.
- [6] Sergey Brin and Lawrence Page. The anatomy of a large-scale hypertextual web search engine. *Computer Networks*, 30(1/7):107–117, 1998.
- [7] M. Burner. Crawling towards eternity: Building an archive of the world wide web. *Web Techniques*, 2(5), 1997.
- [8] Tushar Deepak Chandra and Sam Toueg. Unreliable failure detectors for reliable distributed systems. *Journal of the ACM*, 43(2):225–267, 1996.
- [9] Junghoo Cho and Hector Garcia-Molina. Parallel crawlers. In *Proc. of the 11th International World-Wide Web Conference*, 2002.
- [10] Robert Devine. Design and implementation of DDH: A distributed dynamic hashing algorithm. In David B. Lomet, editor, *Proc. Foundations of Data Organization and Algorithms, 4th International Conference, FODO'93*, volume 730 of *Lecture Notes in Computer Science*, pages 101–114, Chicago, Illinois, USA, 1993. Springer-Verlag.

- [11] Edsger W. Dijkstra. Self-stabilizing systems in spite of distributed control. *Communications of the ACM*, 17(11):643–644, 1974.
- [12] Allan Heydon and Marc Najork. Performance limitations of the Java core libraries. *Concurrency: Practice and Experience*, 12(6):363–373, May 2000.
- [13] Wolfgang Hoschek. The Colt distribution. <http://tilde-hoschek.home.cern.ch/~hoschek/colt/>.
- [14] David Karger, Eric Lehman, Tom Leighton, Matthew Levine, Daniel Lewin, and Rina Panigrahy. Consistent hashing and random trees: Distributed caching protocols for relieving hot spots on the World Wide Web. In *Proc. of the 29th Annual ACM Symposium on Theory of Computing*, pages 654–663, El Paso, Texas, 1997.
- [15] David Karger, Tom Leighton, Danny Lewin, and Alex Sherman. Web caching with consistent hashing. In *Proc. of 8th International World-Wide Web Conference*, Toronto, Canada, 1999.
- [16] Martijn Koster. The Robot Exclusion Standard. <http://www.robotstxt.org/>.
- [17] M. Matsumoto and T. Nishimura. Mersenne twister: A 623-dimensionally equidistributed uniform pseudo-random number generator. *ACMTMCS: ACM Transactions on Modeling and Computer Simulation*, 8:3–30, 1998.
- [18] Marc Najork and Allan Heydon. High-performance web crawling. In J. Abello, P. Pardalos, and M. Resende, editors, *Handbook of Massive Data Sets*. Kluwer Academic Publishers, Inc., 2001.
- [19] Marc Najork and Janet L. Wiener. Breadth-first search crawling yields high-quality pages. In *Proc. of 10th International World Wide Web Conference*, Hong Kong, China, 2001.
- [20] 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, Stanford University, Stanford, CA, USA, 1998.
- [21] Java™ remote method invocation (RMI). <http://java.sun.com/products/jdk/rmi/>.
- [22] Vladislav Shkapenyuk and Torsten Suel. Design and implementation of a high-performance distributed web crawler. In *IEEE International Conference on Data Engineering (ICDE)*, 2002.
- [23] Hongfei Yan, Jianyong Wang, Xiaoming Li, and Lin Guo. Architectural design and evaluation of an efficient Web-crawling system. *The Journal of Systems and Software*, 60(3):185–193, 2002.
- [24] Demetrios Zeinalipour-Yazti and Marios Dikaiakos. Design and implementation of a distributed crawler and filtering processor. In *Proc. of NGITS 2002*, volume 2382 of *Lecture Notes in Computer Science*, pages 58–74, 2002.