

The quest for correct information on the Web: hyper search engines

Massimo Marchiori¹

Department of Pure and Applied Mathematics, University of Padova, Via Belzoni 7, 35131 Padova, Italy

Abstract

Finding the right information in the World Wide Web is becoming a fundamental problem, since the amount of global information that the WWW contains is growing at an incredible rate. In this paper, we present a novel method to extract from a web object its “hyper” informative content, in contrast with current search engines, which only deal with the “textual” informative content. This method is not only valuable per se, but it is shown to be able to considerably increase the precision of current search engines. Moreover, it integrates smoothly with existing search engines technology since it can be implemented on top of every search engine, acting as a post-processor, thus automatically transforming a search engine into its corresponding “hyper” version. We also show how, interestingly, the hyper information can be usefully employed to face the search engines persuasion problem. © 1997 Published by Elsevier Science B.V.

Keywords: World Wide Web; Information retrieval; Search engines; Sep; Browsers and tools; Design principles and techniques; Integration of heterogeneous information sources; Search techniques

1. Introduction

The World Wide Web is growing at phenomenal rates, as witnessed by every recent estimation. This explosion both of Internet hosts and of people using the web, has made crucial the problem of managing such enormous amount of information. As market studies clearly indicate, in order to survive into this informative jungle, web users have to almost exclusively resort on search engines (automatic catalogs of the web) and repositories (human-maintained collections of links usually topics-based). In turn, repositories are now themselves resorting on search engines to keep their databases up-to-date. Thus, the crucial component in the information management is given by search engines.

As all the recent studies on the subject (see e.g. [8])

report, the performance of actual search engines is far from being fully satisfactory. While so far the technological progress has made possible to reasonably deal with the size of the web, in terms of number of objects classified, the big problem is just to properly classify objects in response to the users’ needs: in other words, to give the user the information s/he needs.

In this paper, we argue for a solution to this problem, by means of a new measure of the informative content of a web object, namely its “hyper information”. Informally, the problem with current search engines is that they look at a web object to evaluate, almost just like a piece of text. Even with extremely sophisticated techniques like those already present in some search engine scoring mechanism, this approach suffers from an intrinsic weakness: it doesn’t take into account the *web structure* the object is part of.

The power of the web just relies on its capability of redirecting the information flow via hyperlinks, so

¹ E-mail: max@math.unipd.it

it should appear rather natural that in order to evaluate the informative content of a web object, the web structure has to be carefully analyzed. Instead, so far only very limited forms of analysis of an object *in* the web space have been used, like “visibility” (the number of links pointing to a web object), and the gain in obtaining a more precise information has been very little.

In this paper, instead, we develop a natural notion of “hyper” information which properly takes into account the web structure an object is part of. Informally, the “overall information” of a web object is not composed only by its static “textual information”, but also another “hyper” information has to be added, which is the measure of the potential information of a web object with respect to the web space. Roughly speaking, it measures how much information one can obtain using that page with a browser, and navigating starting from it.

Besides its fundamental characteristic, and future potentials, the hyper information has an immediate practical impact: indeed, it works “on top” of any textual information function, manipulating its “local” scores to obtain the more complete “global” score of a web object. This means that it provides a smooth integration with existing search engines, since it can be used to considerably improve the performance of a search engine simply by post-processing its scores.

Moreover, we show how the hyper information can nicely deal with the big problem of “search engines persuasion” (tuning pages so to cheat a search engine, in order to make it give a higher rank), one of the most alarming factors of degrade of the performance of most search engines.

Finally, we present the result of an extensive testing on the hyper information, used as post-processor in combination with the major search engines, and show how the hyper information is able to *considerably* improve the quality of information provided by search engines.

2. World Wide Web

In general, we consider an (*untimed*) web structure to be a partial function from **URLs**² to se-

quences of bytes. The intuition is that for each URL we can require from the web structure the corresponding object (an **HTML**³ page, a text file, etc.). The function has to be partial because for some URL there is no corresponding object.

In this paper we consider as web structure the World Wide Web structure **WWW**. Note that in general the real **WWW** is a timed structure, since the URL mapping varies with time (i.e. it should be written as WWW_t , for each time instant t). However, we will work under the hypothesis that the **WWW** is *locally time consistent*, i.e. that there is a non-void time interval I such that the probability that $WWW_t(url) = seq$, $t \geq t' \geq t + I \Rightarrow WWW_{t'}(url) = seq$ is extremely high. That is to say, if at a certain time t an URL points to an object seq , then there is a time interval in which this property stays true. Note this doesn't mean that the web structure stays the same, since new web objects can be added.

We will also assume that all the operations resorting on the **WWW** structure that we will employ in this paper, can be performed with extremely high probability within time I : this way, we are acting on a “locally consistent” time interval, and thus it is (with extremely high probability) safe to get rid of the dynamic behavior of the World Wide Web and consider it as an untimed web structure, as we will do.

Note that these assumptions are not restrictive, since empirical observations (cf. [2]) show that **WWW** is indeed locally time consistent, e.g. setting $I = 1$ day (in fact, this is one of the essential reasons for which the World Wide Web works). This, in turn, also implies that the second hypothesis (algorithms working in at most I time) is not restrictive.

A *web object* is a pair (url, seq) , made up by an URL url and a sequence of bytes $seq = WWW(url)$.

In the sequel, we will usually consider understood the **WWW** web structure in all the situations where web objects are considered.

As usual in the literature, when no confusion arises we will sometimes talk of a link meaning the pointed web object (for instance, we may talk about the score of a link meaning the score of the web object it points to).

² <http://www.w3.org/pub/WWW/Addressing/rfc1738.txt>

³ <http://www.w3.org/pub/WWW/MarkUp/>

3. Hyper vs. non-hyper

A great problem with search engines scoring mechanism is that they tend to score *text* more than *hypertext*. Let us explain this better. It is well known that the essential difference between normal texts and hypertext is the relation-flow of information which is possible via hyperlinks. Thus, when evaluating the informative content of a hypertext it should be kept into account this fundamental, *dynamic* behavior of hypertext. Instead, most of search engines tend to simply forget the “hyper” part of a hypertext (the links), by simply scoring its textual components, i.e. providing the *textual information* (henceforth, we say “information” in place of the more appropriate, but longer, “measure of information”).

Note that ranking in a different way words appearing in the **title**⁴, or between **headings**⁵, etc., like all contemporary search engines do, is not in contradiction with what said above, since title, headings and so on are not “hyper”, but they are simply *attributes* of the text.

More formally, the textual information of a web object (*url, seq*) considers only the textual component *seq*, not taking into account the hyper information given by *url* and by *the underlying World Wide Web structure*.

A partial exception is considering the so-called *visibility* (cf. for instance [1]) of a web object. Roughly, the visibility of a web object is the number of other web objects with a link to it. Thus, visibility is in a sense a measure of the importance of a web object in the World Wide Web context. **Excite**⁶ and **WebCrawler**⁷ have been the first search engines to provide higher ranks to web objects with high visibility, now followed by **Lycos**⁸ and **Magellan**⁹.

The problem is that visibility *says nothing* about the *informative content* of a web object. The misleading assumption is that if a web object has a high visibility, then this is a sign of importance and consideration, and so de facto its informative con-

tent must be more valuable than other web objects that have less links pointing to them. This reasoning would be correct if all web objects were known to users and to search engines. But this assumption is clearly false. The fact is that a web object which is not known enough, for instance for its location, is going to have a very low visibility (when it is not completely neglected by search engines), unregarding of its informative content which may be by far better than other web objects.

In a nutshell, visibility is likely to be a synonymous of *popularity*, which is something completely different by *quality*, and thus its usage to give higher score by search engines is a rather poor choice.

3.1. The hyper information

As said, what is really missing in the evaluation of the score of a web object is its hyper part, that is the dynamic information content which is provided by hyperlinks (henceforth, simply links).

We call this kind of information *hyper information*: this information should be added to the *textual information* of the web object, giving its (*overall*) *information* in the World Wide Web. We indicate these three kinds of information as HYPERINFO, TEXTINFO and INFORMATION, respectively. So, for every web object *A* we have that $\text{INFORMATION}(A) = \text{TEXTINFO}(A) + \text{HYPERINFO}(A)$ (note that in general these information functions depend on a specific query, that is to say they measure the informative content of a web object with respect to a certain query: in the sequel, we will always consider such query to be understood).

The presence in a web object of links clearly augments the informative content with the information contained in the pointed web objects (although we have to establish to what extent).

Recursively, links present in the pointed web objects further contribute, and so on. Thus, in principle, the analysis of the informative content of a web object *A* should involve all the web objects that are reachable from it via hyperlinks (i.e., “navigating” in the World Wide Web).

This is clearly unfeasible in practice, so, for practical reasons, we have to stop the analysis at a certain depth, just like in programs for chess analysis we

⁴ <http://www.w3.org/pub/WWW/TR/REC-html32.html#title>

⁵ <http://www.w3.org/pub/WWW/TR/REC-html32.html#headings>

⁶ <http://www.excite.com>

⁷ <http://www.webcrawler.com>

⁸ <http://www.lycos.com>

⁹ <http://www.mckinley.com>

have to stop considering the moves tree after few moves. So, one should fix a certain upper bound for the *depth* of the evaluation. The definition of depth is completely natural: given a web object O , the (relative) depth of another web object O' is the minimum number of links that have to be activated (“clicked”) in order to reach O' from O . So, saying that the evaluation has max depth k means that we consider only the web objects at depth less or equal than k .

By fixing a certain depth, we thus select a suitable finite “local neighborhood” of a web object in the World Wide Web. Now we are faced with the problem of establishing the hyper information of a web object A w.r.t. this neighborhood (the hyper information *with depth* k). We denote this information by with $\text{HYPERINFO}_{[k]}$, and the corresponding overall information (with depth k) with $\text{INFORMATION}_{[k]}$. In most of cases, k will be considered understood, and so we will omit the $[k]$ subscript.

We assume that INFORMATION , TEXTINFO and HYPERINFO are functions from web objects to non-negative real numbers. Their intuitive meaning is that the more information a web object has, the greater is the corresponding number.

Observe that in order to have a feasible implementation, we need that all of these functions are bounded (i.e., there is a number M such that $M \geq \text{INFORMATION}(A)$, for every A). So, we assume without loss of generality that TEXTINFO has an upper bound of 1, and that HYPERINFO is bounded.

3.2. Single links

To start with, consider the simple case where we have only at most one link in every web object.

In the basic case when the depth is one, the most complicated case is when we have one link from the considered web object A to another web object B , i.e.



A naïve approach is just to add the textual information of the pointed web object. This idea, more or less, is tantamount to identify a web object with the web object obtained by replacing the links with the corresponding pointed web objects.

This approach is attracting, but not correct, since it raises a number of problems. For instance, suppose

that A has almost zero textual information, while B has an extremely high overall information. Using the naïve approach, A would have more informative content than B , while it is clear that the user is much more interested in B than in A .

The problem becomes even more evident when we increase the depth of the analysis: consider a situation where B_k is at depth k from A , i.e. to go from A to B one needs k “clicks”, and k is big (e.g. $k > 10$):



Let A, B_1, \dots, B_{k-1} have almost zero informative content, and B_k have very high overall information. Then A would have higher overall information than B_k , which is paradoxical since A has clearly a too weak relation with B_k to be more useful than B_k itself.

The problem essentially is that the textual information pointed by a link cannot be considered as *actual*, since it is *potential*: for the user there is a *cost* to retain the textual information pointed by a link (click and...wait).

The solution to these two factors is: the contribution to the hyper information of a web object at depth k is not simply its textual information, but it is its textual information diminished via a fading factor depending on its depth, i.e. on “how far” is the information for the user (how many clicks s/he has to perform).

Our choice about the law regulating this fading function is that textual information fades exponentially w.r.t. the depth, i.e. the contribution to the hyper information of A given by an object B at depth k is $\mathbf{F}^k \cdot \text{TEXTINFO}(B)$, for a suitable fading factor \mathbf{F} ($0 < \mathbf{F} < 1$).

Thus, in the above example, the hyper information of A is not $\text{TEXTINFO}(B_1) + \dots + \text{TEXTINFO}(B_k)$, but $\mathbf{F} \cdot \text{TEXTINFO}(B_1) + \mathbf{F}^2 \cdot \text{TEXTINFO}(B_2) + \dots + \mathbf{F}^k \cdot \text{TEXTINFO}(B_k)$, so that the overall information of A is not necessarily greater than that of B_k .

As an aside, note that when calculating the overall information of a web object A , its textual information can be nicely seen as a special degenerate case of hyper information, since $\text{TEXTINFO}(A) = \mathbf{F}^0 \cdot \text{TEXTINFO}(A)$ (viz., the object is at “zero distance” from itself).

3.3. More motivations

There is still one point to clarify. We have introduced an exponentially fading law, and seen how it behaves well with respect to our expectations. But one could ask why the fading function has to be exponential, and not of different form, for instance polynomial (e.g., multiplying by F/k instead that by F^k). Let us consider the overall information of a web object with depth k : if we have



it is completely reasonable to assume that the overall information of A (with depth k) is given by the textual content of A plus the faded overall information of B (with depth $k - 1$). It is also reasonable to assume that this fading can be approximated by multiplying by an appropriate fading constant F ($0 < F < 1$). So, we have the recursive relation

$$\text{INFORMATION}_{[k]}(A) = F \cdot \text{INFORMATION}_{[k-1]}(B)$$

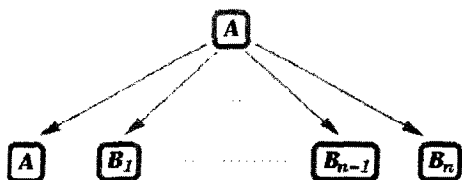
Since readily for every object A we have $\text{INFORMATION}_{[0]}(A) = \text{TEXTINFO}(A)$, eliminating recursion we obtain that if



then $\text{INFORMATION}(A) = \text{TEXTINFO}(A) + F \cdot (\text{TEXTINFO}(B_1) + F \cdot (\text{TEXTINFO}(B_2) + F \cdot (\dots \text{TEXTINFO}(B_k)))) = \text{TEXTINFO}(A) + F \cdot \text{TEXTINFO}(B_1) + F^2 \cdot \text{TEXTINFO}(B_2) + \dots + F^k \cdot \text{TEXTINFO}(B_k)$, which is just the exponential fading law that we have adopted.

3.4. Multiple links

Now we turn to the case where there is more than one link in the same web object. For simplicity, we assume the depth is one. So, suppose one has the following situation



What is the hyper information in this case? The easiest answer, just sum the contribution of every

link (i.e. $F \cdot \text{TEXTINFO}(B_1) + \dots + F \cdot \text{TEXTINFO}(B_n)$), isn't feasible since we want the hyper information to be *bounded*.

This would seem in contradiction with the interpretation of a link as potential information that we have given earlier: if one has many links that can be activated, then one has all of their potential information. However, this paradox is only apparent: the user cannot get all the links at the same time, but has to *sequentially select* them. In other words, *nondeterminism has a cost*. So, in the best case the user will select the most informative link, and then the second more informative one, and so on. Suppose for example that the more informative link is B_1 , the second one is B_2 and so on (i.e., we have $\text{TEXTINFO}(B_1) \geq \text{TEXTINFO}(B_2) \geq \dots \geq \text{TEXTINFO}(B_n)$). Thus, the hyper information is $F \cdot \text{TEXTINFO}(B_1)$ (the user selects the best link) plus $F^2 \cdot \text{TEXTINFO}(B_2)$ (the second time, the user selects the second best link) and so on, that is to say

$$F \cdot \text{TEXTINFO}(B_1) + \dots + F^n \cdot \text{TEXTINFO}(B_n)$$

Nicely, evaluating the score this way gives a bounded function, since for any number of links, the sum cannot be greater than $F/(F + 1)$.

Note that we chose the best sequence of selections, since hyper information is the best "potential" information, so we have to assume the user does the best choices: we cannot use e.g. a random selection of the links, or even other functions like the average between the contributions of the each link, since we cannot impose that every link has to be relevant. For instance, if we did so, accessory links with zero score (e.g. think of the "powered with Netscape"¹⁰-links) would devalue by far the hyper information even in presence of highly scored links, while those accessory links should simply be ignored (as the above method, consistently, does).

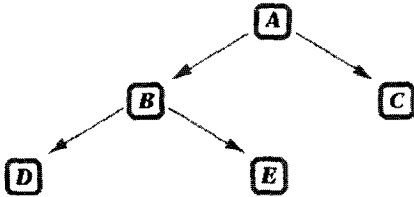
3.5. The general case

The general case (multiple links in the same page and arbitrary depth k) is treated accordingly to what previously seen. Informally, all the web objects at depth less of equal than k are ordered w.r.t. a

¹⁰ http://www.netscape.com/comprod/products/navigator/version_3.0/images/netnow3.gif

“sequence of selections” such that the corresponding hyper information is the highest one.

For example, consider the following situation:



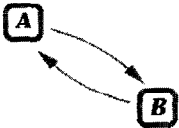
with $F = 0.5$, $\text{TEXTINFO}(B) = 0.4$, $\text{TEXTINFO}(C) = 0.3$, $\text{TEXTINFO}(D) = 0.2$, $\text{TEXTINFO}(E) = 0.6$.

Then via a “sequence of selections” we can go from A to B , to C , to E and then to D , and this is readily the best sequence that maximizes the hyper information, which is $0.5 \cdot \text{TEXTINFO}(B) + 0.52 \cdot \text{TEXTINFO}(C) + 0.53 \cdot \text{TEXTINFO}(E) + 0.54 \cdot \text{TEXTINFO}(D) (= 0.3625)$.

4. Refining the model

There is a number of subtleties that for clarity of exposition we haven’t considered so far.

The first big problem comes from possible duplications. For instance, consider the case of *backwards* links: suppose that a user has set up two pages A and B , with a link in from A to B , and then a “come back” link in B from B to A :



This means that we have a kind of recursion, and that the textual information of A is repeatedly added (although with increasingly higher fading) when calculating the hyper information of A (this because from A we can navigate to B and then to A and so on...). Another example is given by *duplicating links* (e.g. two links in a web object pointing to the same web object). The solution to avoid these kinds of problems, is not to consider all the “sequence of selections”, but only those *without repetitions*. This is consistent with the intuition of hyper information to measure in a sense how much a web object is far from a page: if one has already reached an object, s/he already has got its information, and so it makes no sense to get it again.

Another important issue is that the same definition of link in a web object is fairly from trivial. A link present in a web object O is said to be *active* if the web objects it points to can be accessed by viewing (url, seq) with an **HTML**¹¹ browser (e.g., **Netscape**¹², **Navigator**¹³ or **Microsoft Internet Explorer**¹⁵). This means, informally, that once we view O with the browser, we can activate the link by clicking over it. The previous definition is rather operational, but it is much more intuitive than a formal technical definition which can be given by tediously specifying all the possible cases according to the **HTML specification**¹⁶ (note a problem complicating a formal analysis is that one cannot assume that seq is composed by legal HTML code, since browsers are error-tolerating).

Thus, the links mentioned in the paper should be only the active ones.

Also, there are many different kinds of links, and each of them would require a specific treatment. For instance,

- *Local* links (links pointing to some point in the same web object, using the **#-specifier**¹⁷) should be readily ignored.
- *Frame*¹⁸ links should be automatically expanded, i.e. if A has a frame link to B , then this link should be replaced with a proper expansion of B inside A (since a frame link is automatically activated, its pointed web object is just part of the original web object, and the user does not see any link at all).
- Other links that are automatically activated are for instance the *image*¹⁹ links, i.e. source links of image tags, the *background*²⁰ links and so on: they should be treated analogously to frame links. However, since usually it is a very hard problem to recover some useful information from images

¹¹ <http://www.w3.org/pub/WWW/Markup/>

¹² <http://www.netscape.com>

¹³ http://www.netscape.com/comprod/products/navigator/version_3.0/index.html

¹⁴ <http://www.microsoft.com>

¹⁵ <http://www.microsoft.com/ie>

¹⁶ <http://www.w3.org/pub/WWW/Markup/>

¹⁷ <http://www.w3.org/pub/WWW/Addressing/rfc1738.txt>

¹⁸ http://home.netscape.com/assist/net_sites/frames.html

¹⁹ <http://www.w3.org/pub/WWW/TR/REC-html32.html#img>

²⁰ <http://www.w3.org/pub/WWW/TR/REC-html32.html#body>

(cf. [6]), in a practical implementation of the hyper information all these links can be ignored (note this doesn't mean that the textual information function has to ignore such links, since it can e.g. take into account the name of the pictures). This also happens for active links pointing to images, movies, sounds etc., which can be syntactically identified by their file extensions .gif, .jpg, .tiff, .avi, .wav and so on.

4.1. Search engines persuasion

A big problem that search engines have to face is the phenomenon of so-called *sep* (search engines persuasion). Indeed, search engines have become so important in the advertisement market that it has become essential for companies to have their pages listed in top positions of search engines, in order to get a significant web-based promotion. Starting with the already pioneering work of Rhodes [5], this phenomenon is now boosting at such a rate to have provoked serious problems to search engines, and has revolutioned the web design companies, which are now specifically asked not only to design good web sites, but also to make them rank high in search engines. A vast number of new companies was born just to make customer web pages as visible as possible. More and more companies, like **Exploit**²¹, **Allwilk**²², **Northern Webs**²³, **Ryley & Associates**²⁴, etc., explicitly study ways to rank high a page in search engines. **OpenText**²⁵ arrived even to sell “preferred listings”, i.e. assuring a particular entry to stay in the top ten for some time, a policy that has provoked some controversies (cf. [9]).

This has led to a bad performance degradation of search engines, since an increasingly high number of pages is designed to have an artificially high textual content. The phenomenon is so serious that search engines like **InfoSeek**²⁶ and **Lycos**²⁷ have introduced penalties to face the most common of this

persuasion techniques, “spamdexing” [4,7,8], i.e. the artificial repetition of relevant keywords. Despite these efforts, the situation is getting worst and worst, since more sophisticated techniques have been developed, which analyze the behavior of search engines and tune the pages accordingly. This has also led to the situation where search engines maintainers tend to assign penalties to pages that rank “too high”, and at the same time to provide less and less details on their information functions just in order to prevent this kind of tuning, thus in many cases penalizing pages of high informative value that were not designed with “persuasion” intentions. For a comprehensive study of the sep phenomenon, we refer to [3].

The hyper methodology that we have developed is able to a certain extent to nicely cope with the sep problem. Indeed, maintainers can keep details of their TEXTINFO function hidden, and make public the information that they are using a HYPERINFO function.

The only precaution is to distinguish between two fundamental types of links, assigning different fading factors to each of them. Suppose to have a web object (*url, seq*). A link contained in *seq* is called *outer* if it has not the same **domain**²⁸ of *url*, and *inner* in the other case. That is to say, inner links of a web objects point to web objects in the same site (its “local world”, so to say), while outer links point to web objects of other sites (the “outer world”).

Now, inner links are dangerous from the sep point of view, since they are on the direct control of the site maintainer. For instance, a user that wants to artificially increase the hyper information of a web object *A* could set up a very similar web object *B* (i.e. such that $\text{TEXTINFO}(A) \approx \text{TEXTINFO}(B)$), and put a link from *A* to *B*: this would increase the score of *A* by roughly $F \cdot \text{TEXTINFO}(A)$.

On the other hand, outer links do not present this problem since they are out of direct control and manipulation (at least in principle, cf. [2]).

Thus, one should consequently assign a *very low* or even *null* fading factor (F_{in}) to the inner links, and a reasonably high fading factor (F_{out}) to the outer links. Indeed, in our practical experimentations we saw that setting F_{in} to zero, i.e. completely omitting

²¹ <http://www.exploit.com>

²² <http://www.allwilk.com>

²³ <http://www.digital-cafe.com/~webmaster/norweb01.htm>

²⁴ <http://www.ryley.com>

²⁵ <http://www.opentext.com>

²⁶ <http://www.infoseek.com>

²⁷ <http://www.lycos.com>

²⁸ <http://www.dns.net/dnsrd/>

the inner link contributions, gave very good results (although the “best” value was according to our test was near to 0.1). Setting F_{in} to zero also gave the advantage of making the implementation of the hyper information quite faster, since most of the links in web objects are inner. As far as outer links are concerned, we set F_{out} to 0.75. Again, from our tests it resulted that similar values did not significantly affect the bounty of the hyper information.

We said earlier that the information about the use of the hyper information could be given as white box to the external users (while keeping as black box the details of the TEXTINFO function). This way, search engines persuasion has the effect of *reshaping the web*, by considerably improving its connectivity. Indeed, as noticed in [1], at present the inter-connectivity is rather poor, since almost 80% of sites contain no outer link (!), and a relatively small number of web sites is carrying most of the load of hypertext navigation.

Using hyper information thus forces the sites that want to rank better to *improve their connectivity*, improving the overall web structure.

5. Testing

The hyper information has been implemented as *post-processor* of the main search engines now available, i.e. remotely querying them, extracting the corresponding scores (i.e., their TEXTINFO function), and calculating the hyper information and therefore the overall information.

Indeed, as said in the introduction, one of the most appealing features of the hyper methodology is that it can be implemented “on top” of existing scoring functions. Note that, strictly speaking, scoring functions employing visibility, like those of **Excite**²⁹, **WebCrawler**³⁰ and **Lycos**³¹, are not pure “textual information”. However, this is of little importance: although visibility, as shown, is not a good choice, it provides information which is disjoint to the hyper information, and so we can view such scoring functions like providing purely textual information

slightly perturbed (improved?) with some other kind of WWW-based information.

The search engines for which a post-processor was developed were: **Excite**³², **HotBot**³³, **Lycos**³⁴, **WebCrawler**³⁵, and **OpenText**³⁶. This includes all of the major search engines, but for **AltaVista**³⁷ and **InfoSeek**³⁸, which unfortunately do not give the user access to the scores, and thus cannot be remotely post-processed. The implemented model included all the accessory refinements seen in Section 4.

We then conducted some tests in order to see how the values of F_{in} and F_{out} affected the quality of the hyper information. The depth and fading parameters can be flexibly employed to tune the hyper information. However, it is important for its practical use that the behaviour of the hyper information is tested when the depth and the fading factors are fixed in advance.

We randomly selected 25 queries, collected all the corresponding rankings from the aforementioned search engines, and then run several tests in order to maximize the effectiveness of the hyper information. We arrived to select the following global parameters for the hyper information: $F_{in} = 0$, $F_{out} = 0.75$, and depth one. Although the best values were slightly different ($F_{in} = 0.1$, depth two), the differences were almost insignificant. On the other hand, choosing $F_{in} = 0$ and depth one had great advantages in terms of execution speed of the hyper information.

After this setting, we chose other 25 queries, and tested again the bounty of the hyper information with these fixed parameters: the results showed that these initial settings also gave extremely good results for these new set of queries.

While our personal tests clearly indicated a considerable increasing of precision for the search engines, there was obviously the need to get external confirmation of the effectiveness of the approach. Also, although the tests indicated that slightly perturbing the fading factors did not substantially affect the hyper information, the bounty of the values of the

²⁹ <http://www.excite.com>

³⁰ <http://www.webcrawler.com>

³¹ <http://www.lycos.com>

³² <http://www.excite.com>

³³ <http://www.hotbot.com>

³⁴ <http://www.lycos.com>

³⁵ <http://www.webcrawler.com>

³⁶ <http://www.opentext.com>

³⁷ <http://www.altavista.com>

³⁸ <http://www.infoseek.com>

fading factors could somehow have been dependent on our way of choosing queries.

Therefore, we performed a deeper and more challenging test. We considered a group of thirty persons, and asked each of them to arbitrarily select five different topics to search on in the World Wide Web.

Then each person had to perform the following test: s/he had to search for relevant information on each chosen topic by inputting suitable queries to our post-processor, and give an evaluation mark to the obtained ranking lists. The evaluation consisted in an integer ranging from 0 (terrible ranking, of no use), to 100 (perfect ranking).

The post-processor was implemented as a cgi-script: the prompt asks for a query, and returns two ranking lists relative to that query, one per column in the same page (**frames**³⁹ are used). In one column there is the original top ten of the search engine. In the other column, the top ten of the ranking obtained by: (1) taking the first 100 items of the ranking of the search engine, and (2) post-processing them with the hyper information.

We made serious efforts in order to provide each user with the most natural conditions in order to evaluate the rankings. Indeed, to avoid “psychological pressure” we provided each user with a password, so that each user could remotely perform the test on his/her favorite terminal, in every time period (the test process could be frozen via a special button, terminating the session, and resumed whenever wanted), with all the time necessary to evaluate a ranking (e.g. possibly via looking at each link using the browser). Besides the evident beneficial effects on the bounty of the data, this also had the side-effect that we hadn’t to personally take care of each test, something which could have been extremely time consuming.

Another *extremely important* point is that the post-processor was implemented to make the test *completely blind*. This was achieved in the following way.

Each time a query was run, the order in which the five search engines had to be applied was randomly chosen. Then, for each of the five search engines the following process occurred:

- The data obtained by the search engine were “filtered” and presented in a standard “neutral” format. The same format was adopted for the

post-processed data, so that the user couldn’t see any difference between the two rankings but for their links content.

- The columns where to view the original search engine results and the post-processed results were randomly chosen.
- The users were provided with the only information that we were performing a test on evaluating several different scoring mechanisms, and that the choice of the column where to view each ranking was random. This also avoided possible “chain” effects, where a user, after having given for a certain number of times higher scores to a fixed column, can be unconsciously led to give it bonuses even if it actually doesn’t deserve them.

Fig. 1 shows an example snapshot of a session, where the selected query is “search engine score”: the randomly selected search engine in this case was **HotBot**⁴⁰, and its original ranking (filtered) is in the left column.

The final results of the test are pictured in Fig. 2. As it can be seen, the chart clearly shows (especially in view of the blind evaluation process), that the hyper information is able to considerably improve the evaluation of the informative content.

The precise data regarding the evaluation test are shown in the next table:

	Excite	HotBot	Lycos	WebCrawler	OpenText	Average
Normal	80.1	62.2	59.0	54.2	63.4	63.2
Hyper	85.2	77.3	75.4	68.5	77.1	76.7

The next table shows the evaluation increment for each search engine w.r.t. its hyper version, and the corresponding standard deviation:

	Excite	HotBot	Lycos	WebCrawler	OpenText	Average
Evaluation increment	+5.1	+15.1	+16.4	+14.3	+13.7	+12.9
Standard deviation	2.2	4.1	3.6	1.6	3.0	2.9

As it can be seen, the small standard deviations are a further empirical evidence of the superiority of the hyper search engines over their non-hyper versions.

³⁹ http://home.netscape.com/assist/net_sites/frames.html

⁴⁰ <http://www.hotbot.com>

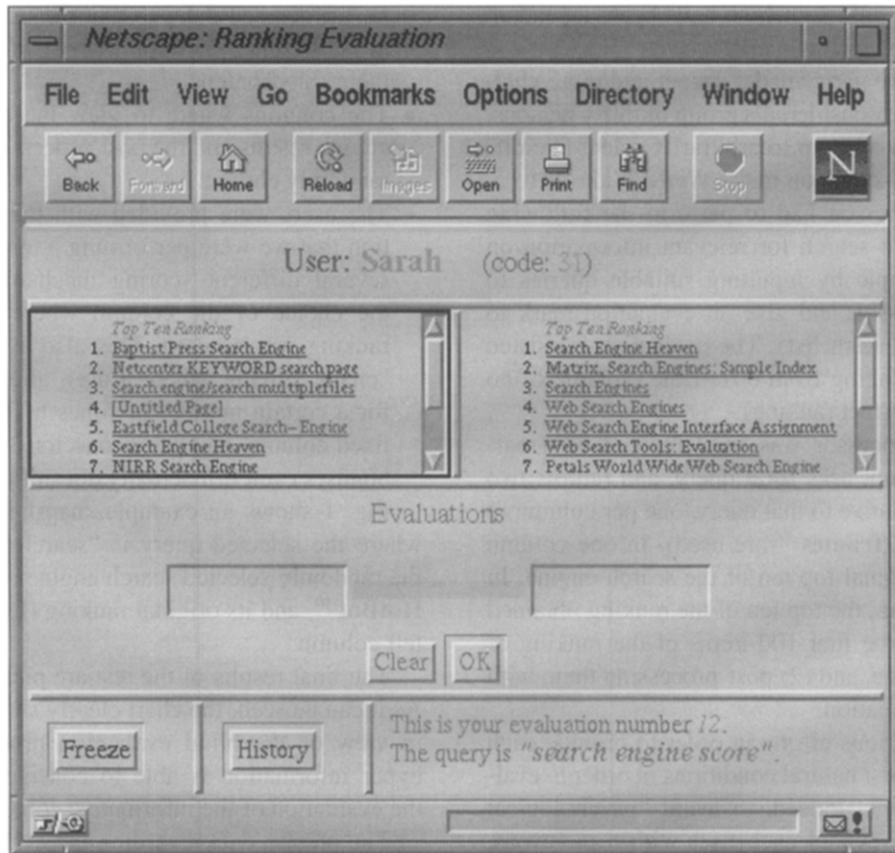


Fig. 1. Snapshot from a test session.

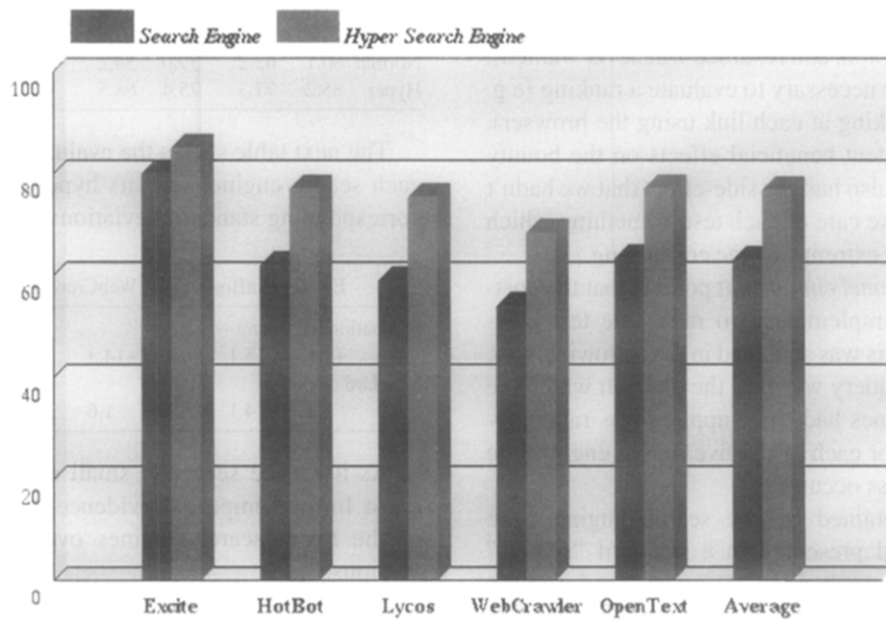


Fig. 2. Evaluation of search engines vs. their hyper versions.

6. Conclusion

We have presented a new way to considerably improve the quality of search engines evaluations by considering the “hyper information” of web objects. We have shown that besides its practical success, the method presents several other interesting features:

- It allows to focus separately on the “textual” and “hyper” components.
- It allows to improve current search engines in a smooth way, since it works “on top” of any existing scoring function.
- It can be used for locally improving performances of search engines, just using a local post-processor on the client side (like done in this paper)
- Analogously, it can be used to improve so-called “meta searchers”, like **SavySearch**⁴¹, **Meta-Crawler**⁴², **WebCompass**⁴³, **NlightN**⁴⁴ and so on. These suffer even more from the problem of good scoring, since they have to mix data coming from different search engines, and the result is in most of cases not very readable. The method can considerably improve the effectiveness and utility for the user of such meta-engines.
- Last but not least, it behaves extremely well with respect to search engines persuasion attacks.

References

- [1] T. Bray (mailto:tbray@textuality.com), *Measuring the Web* (http://www5conf.inria.fr/fich_html/papers/P9/Overview.html), *5th International World Wide Web Conference* (<http://www.w3.org/pub/Conferences/WWW5/>), May, Paris, 1996
- [2] M. Marchiori (mailto:max@math.unipd.it), *The hyper information: theory and practice*, Technical Report No. 46, Dept. of Pure & Applied Mathematics, University of Padova, 1996
- [3] M. Marchiori (mailto:max@math.unipd.it), *Security of World Wide Web search engines*, *3rd International Conference on Reliability, Quality and Safety of Software-Intensive Systems*, Athens, Greece, Chapman & Hall (<http://www.thomson.com:8866/chaphall/default.html>), 1997
- [4] K. Murphy (mailto:kathleen@webweek.com), *Cheaters never win* (<http://www.webweek.com/96May20/undercon/heaters.html>), *Web Week* (<http://www.webweek.com>), May, 1996
- [5] J. Rhodes (mailto:deadlock@deadlock.com), *The Art of Business Web Site Promotion* (<http://deadlock.com/~deadlock/promote>), *Deadlock Design* (<http://deadlock.com/~deadlock>), 1996.
- [6] S. Sclaroff (mailto:sclaroff@cs.bu.edu), *World Wide Web image search engines* (<http://www.cs.bu.edu/techreports/95-016-www-image-search-engines.ps.z>), *NSF Workshop on Visual Information Management*, Cambridge, MA, 1995
- [7] D. Sullivan (mailto:danny@calafia.com), *The Webmaster's Guide to Search Engines and Directories* (<http://calafia.com/webmasters/>), Calafia Consulting (<http://calafia.com>), 1996
- [8] G. Venditto (mailto:venditto@iw.com), *Search engine showdown* (<http://www.internetworld.com/1996/05/showdown.html>), *Internet World* (<http://www.internetworld.com>), Vol. 7(5) (<http://www.internetworld.com/1996/05/toc.html>), May, 1996.
- [9] N. Wingfield (mailto:nickw@cnet.com), *Engine sells results, draws fire* (<http://www.news.com/News/Item/0,4,1635,00.html>), C|net Inc. (<http://www.news.com>), June, 1996.



Massimo Marchiori received the M.S. in Mathematics with Highest Honors, and the Ph.D. in Computer Science, both from the University of Padua. His research interests include World Wide Web and intranets (information retrieval, search engines, metadata, digital libraries), visualisation, programming languages (constraint programming, visual programming, functional programming, logic programming), genetic algorithms, and rewriting systems. He has published over twenty refereed papers on the above topics in various journals and proceedings of international conferences.

⁴¹ <http://www.cs.colostate.edu/~dreiling/smartform.html>

⁴² <http://www.metacrawler.com>

⁴³ http://arachnid.qdeck.com/qdeck/demosoft/webcompass_lite/

⁴⁴ <http://www.nlightn.com>