

PicASHOW: Pictorial Authority Search by Hyperlinks On the Web

RONNY LEMPEL

Department of Computer Science, The Technion, Haifa

and

AYA SOFFER

IBM Research Lab in Haifa

We describe PicASHOW, a fully automated WWW image retrieval system that is based on several link-structure analyzing algorithms. Our basic premise is that a page p displays (or links to) an image when the author of p considers the image to be of value to the viewers of the page. We thus extend some well known link-based WWW *page retrieval* schemes to the context of image retrieval.

PicASHOW's analysis of the link structure enables it to retrieve relevant images even when those are stored in files with meaningless names. The same analysis also allows it to identify *image containers* and *image hubs*. We define these as Web pages that are rich in relevant images, or from which many images are readily accessible.

PicASHOW requires no image analysis whatsoever and no creation of taxonomies for pre-classification of the Web's images. It can be implemented by standard WWW search engines with reasonable overhead, in terms of both computations and storage, and with no change to user query formats. It can thus be used to easily add image retrieving capabilities to standard search engines.

Our results demonstrate that PicASHOW, while relying almost exclusively on link analysis, compares well with dedicated WWW image retrieval systems. We conclude that link analysis, a proven effective technique for Web page search, can improve the performance of Web image retrieval, as well as extend its definition to include the retrieval of image hubs and containers.

Categories and Subject Descriptors: H.2.8 [Database Management]: Database Applications—*Image databases*; H.3.1 [Information Storage and Retrieval]: Content Analysis and Indexing—*Indexing methods*; H.3.3 [Information Storage and Retrieval]: Information Search and Retrieval—*Search process*; H.3.5 [Information Storage and Retrieval]: Online Information Services—*Web-based services*; H.5.4 [Information Interfaces and Presentation]: Hypertext/Hypermedia—*Theory*

General Terms: Algorithms

Additional Key Words and Phrases: Image retrieval, link structure analysis, hubs and authorities, image hubs

Authors' addresses: R. Lempel, Department of Computer Science, The Technion, Haifa 32000, Israel; email: rlempel@cs.technion.ac.il; A. Soffer, IBM Research Lab in Haifa, Matam Park, Haifa 31905, Israel; email: ayas@il.ibm.com.

Permission to make digital/hard copy of part or all of this work for personal or classroom use is granted without fee provided that the copies are not made or distributed for profit or commercial advantage, the copyright notice, the title of the publication, and its date appear, and notice is given that copying is by permission of the ACM, Inc. To copy otherwise, to republish, to post on servers, or to redistribute to lists, requires prior specific permission and/or a fee.

© 2002 ACM 1046-8188/02/0100-0001 \$5.00

1. INTRODUCTION

The WWW is host to millions of images on almost every conceivable topic. Finding effective methods to retrieve these images has attracted many research efforts over the past few years. Such research has led to academic image retrieval systems (Columbia University's webSEEK system¹), to general search engines with image retrieval capabilities (AltaVista², Lycos Multimedia³), and to search engines dedicated to WWW image retrieval (Scour⁴, Ditto⁵). There are three main approaches for WWW image search and retrieval:

- (1) *Text-based retrieval*. This approach annotates images with text derived from the HTML documents which contain (display) them, and then applies text-based retrieval algorithms to the annotated collection of images. The derived text can include the caption of the image, text surrounding the image, the entire text of the containing page, the filename of the containing HTML document and the filename of the image itself.
- (2) *Content-based image retrieval (CBIR)*, [Flickner et al. 1995; Pentland et al. 1996; Rui et al. 1999]). This approach applies image analysis techniques in order to extract visual features (e.g., color, texture, orientation, shape) from the images. The features are extracted in a preprocessing stage, and stored in the retrieval system's database. The extracted features (e.g., the color histogram of the image) are usually of high dimensionality, and in order to allow scalability of these systems (in terms of storage space and query processing times), some sort of dimension reduction is usually performed on the data (e.g., Kanth et al. [1998]).
- (3) *Manually annotated image collections*. There are several firms that specialize in providing visual content to a diverse range of image consumers. The two largest are Getty Images⁶ and Corbis.⁷ These firms maintain archives where images are indexed and retrieved by keywords, which are manually assigned to each image. While end users may search the image archives of these firms (through the firms' Web sites), the main customers are companies and professionals who require high volumes of diverse images. These corporate customers may, in turn, offer Web users image retrieval services, presenting the visual content provided to them by the image archiving firms.

One implication of the difference between the retrieval approaches is their support of different types of queries. Text-based retrieval systems, as well as the commercial image providers, support natural, topic-descriptive queries. These queries are friendly and familiar to the typical surfer of the Web. On the other hand, CBIR supports either queries which are formulated in terms

¹<http://www.ctr.columbia.edu/webseek/>

²AltaVista Company, <http://www.altavista.com/>

³Lycos Inc. <http://multimedia.lycos.com/>

⁴Scour Inc. <http://www.scour.com/>

⁵ditto.com visual search engine, <http://www.ditto.com/>

⁶Getty Images Inc. <http://www.getty-images.com/>

⁷Corbis Inc. <http://www.corbis.com/>

of the extracted visual features, or similarity queries, in which a sample image is presented and the system is required to retrieve images with similar visual features.

Web Image Search. PicToSeek [Gevers and Smeulders 1999] is an example of a pure content-based image retrieval system. It classifies images into portraits, photographs of indoor/outdoor scenes, and synthetic images. It extracts many visual features from the images, and supports queries by image-examples and by image features.

Many WWW image search engines combine text-based retrieval and CBIR into an integrated system. In *webSeek* [Smith and Chang 1996], for example, each image is processed and its visual features are extracted. Each image is also associated with the text in its containing page. The images are then classified into topics from a taxonomy that was developed for this purpose. The *WebSeer* system [Frankel et al. 1996] uses associated text and feature extraction to support complex queries, which state both the search topic and some visual properties of the desired images. The system depicted in Cascia et al. [1988] unifies the textual representations and visual representations of images into a single representative vector. The textual representations are derived from latent semantic analysis of the text in the containing HTML pages, while the visual representations are dimensionally reduced color and orientation histograms. This unified representation enables the utilization of possible statistical couplings between the textual contents of the containing pages and the visual properties of the images.

Harmandas et al. [1997] suggest a text-based image retrieval system in which connectivity information is used to induce textual annotations of images. Each image i in their approach is assigned a weighted vector of representing terms, which is some function of the combined text of all of the pages that contain i and of the pages that link to pages containing i . While this scheme does consider hyperlinks, it is essentially a text-based retrieval scheme in which hyperlinks are used to induce textual annotations, without any analysis being done on the link-structure per se.

Link Structure Analysis in Web Page Search. Recent work in Web search has demonstrated that link structure analysis is very effective in finding authoritative Web pages. Information such as which pages are linked to others is commonly used to augment search algorithms, and has significantly improved the ability of search engines to rank quality pages at the top of their search results. Link-structure analysis is based on the notion that a link from page p to page q can be viewed as an endorsement of q by p , and as some form of positive judgment by p of q 's content.

Two important types of techniques in link-structure analysis are co-citation based schemes, and random-walk based schemes. The main idea behind co-citation based schemes is the notion that when two pages p_1 and p_2 both point to some page q , it is reasonable to assume that p_1 and p_2 share a mutual topic of interest. Likewise, when p links to both q_1 and q_2 , it is probable that q_1 and q_2 share some mutual topic. An important work in the context of co-citation based schemes was Jon Kleinberg's introduction of the notions of *hubs* and

authorities [Kleinberg 1999] as two distinct types of Web pages. Authorities, or authoritative pages, are Web pages that contain high-quality information regarding some topic. Hubs, on the other hand, may not directly contain information but are rather resource lists, linking to authorities on a topic without necessarily displaying the information itself. Kleinberg devised an algorithm aimed at finding authoritative pages, and researchers from IBM's Almaden Research Center have implemented Kleinberg's algorithm in various projects, most notably CLEVER.⁸

Random walk based schemes model the Web (or part of it) as a graph (where pages are nodes and links are edges), and apply some random walk model to the graph. Pages are then ranked by the probability of visiting them in the modeled random walk. The most notable algorithm of this type is PageRank [Brin and Page 1998], which is an important part of the ranking function and of the success of the Google search engine.⁹ Both Kleinberg's algorithm and PageRank are described in detail in Section 2.

Co-citation reasoning was combined with random walk theory in SALSA [Lempel and Moran 2000], to separate the random walk based rankings of hubs and authorities. Rafiei and Mendelzon [2000] have also integrated co-citation and random walks in their work on computing page reputations.

Our Approach: Link Structure Analysis in Web Image Search. In this paper we present PicASHOW, a pictorial retrieval system that searches for images (pictures) on the Web using hyperlink-structure analysis. PicASHOW applies co-citation based approaches and PageRank influenced methods. Our basic premise is that a page p displays (or links to) an image when the author of p considers the image to be of value to the viewers of the page. We further contend that the standard reasoning behind the co-citation measure applies to images just as it does to HTML pages:

- Images which are co-contained in pages are likely to be related to the same topic.
- Images which are contained in pages that are co-cited by a certain page are likely related to the same topic.

In addition, in the spirit of PageRank, we assume that images which are contained in authoritative pages on topic t are good candidates to be quality images on that topic.

In the next sections, we describe several link-structure based WWW image retrieval schemes. Following are the highlights of the PicASHOW approach:

- Our method can be implemented, with reasonable overhead, by standard WWW search engines. It can thus be used to add image retrieving capabilities to these engines. We elaborate on this in Section 3.2.
- Our schemes require no image analysis whatsoever. This eliminates the need to deal with high-dimension image descriptions, and with the complexity

⁸<http://www.almaden.ibm.com/cs/k53/clever.html>

⁹Google Inc. <http://www.google.com/>

which such representations introduce in terms of memory requirements, pre-processing overhead, query processing and retrieval operations.

- No change to the query format is required. The same queries which are used to retrieve pages, will be used to retrieve images. In particular, users do not need to present the system with sample images, nor do they need to formulate queries in terms of image properties.
- There is no need to create taxonomies for preclassification of the wealth of images on the Web.
- We do not rely solely on file names and image captions assigned by content creators. Thus, we are able to find images related to a query with meaningless file names such as “myimages/image1” (most text-based image search engines will miss these). We can also find images with titles that are only semantically related to the query. For example images labeled “Bridal Veil Falls”, when searching for images of Yosemite.
- In addition to finding authoritative images, we are also able to locate image containers and image hubs. We define these as Web *pages* that are rich in relevant images, or from which many images are readily accessible. See Section 4 for more details.
- A natural modification of our methods allows for the support of similarity queries,¹⁰ where users present PicASHOW with URLs of images on the topic in question. The system will then find other authoritative images on the same topic. We believe this is a very useful feature. Section 6 elaborates on the details.

The remainder of this paper is organized as follows. In Section 2 we provide some background on link-structure analysis when searching for Web pages. In Section 3 we formally define the image collections that are to be analyzed, explain how such collections are assembled from a given query, and present our proposed image ranking schemes. Section 4 introduces the concept of image hubs and image containers and describes how we identify such hubs and containers. Section 5 brings an informal comparative evaluation of the image ranking schemes (see below). In Section 6 we discuss the pros and cons of our method, and suggest interesting extensions of this research direction. Appendix A lists the URLs of the images which are displayed in this paper.

We do not provide any formal evaluation section in this paper since there are no benchmarks for testing Web-based image retrieval systems, thus most evaluations are qualitative. Rather, we provide three kinds of results:

- (1) Section 4 brings the URLs of several image hubs and containers found by PicASHOW.
- (2) In Section 5 we report precision@10 values of our image ranking schemes on a diverse range of queries, comparing the effectiveness of those schemes.

¹⁰Note that this is different from the similarity queries mentioned in the context of CBIR systems—our method will find images on the same topic as the sample images rather than images with similar visual properties.

- (3) Appendix A shows the URLs of images that were retrieved by several sample queries. For comparison purposes, we also list the URLs of the results of some of these queries on commercial Web search engines.

The images whose URLs are given in Appendix A and several images from the image containers mentioned in Section 4 were showcased in an earlier version of this paper [Lempel and Soffer 2001].¹¹ Unfortunately, legal concerns have prompted the publishers of this journal to ask us to remove the actual images from the manuscript. While we do provide the URLs of the images, many of these URLs may very shortly become stale because of the volatile nature of the ever-changing Web. We feel that this paper's results are best appreciated with the retrieved images at hand, and encourage the readers to obtain a copy of Lempel and Soffer [2001].

2. LINK ANALYSIS FOR FINDING AUTHORITATIVE WEB PAGES

This section provides some technical background on applications of WWW link-structure analysis when searching for Web pages. Specifically, we provide a brief overview of two link-structure analyzing approaches: PageRank [Brin and Page 1998] and Kleinberg's Mutual Reinforcement approach [Kleinberg 1999]. This background is required in order to describe our image ranking schemes which are inspired by these approaches. Indeed, for each of the two approaches described, we also describe the main points which we adapted and evolved in our image ranking schemes.

2.1 PageRank

PageRank [Brin and Page 1998] is an important part of the ranking function of the Google search engine. The PageRank of a page p is the probability of visiting p in a random walk of the entire Web, where the set of states of the random walk is the set of pages, and each random step is of one of the following two types:

- (1) From the given state s , choose at random an outgoing link of s and follow that link to the destination page.
- (2) Choose a Web page uniformly at random, and jump to it.

PageRank chooses a parameter d , $0 < d < 1$, and each state transition is of the first transition type with probability d and of the second type with probability $1 - d$. The PageRanks obey the following formula (where page p has incoming links from pages q_1, \dots, q_k and N is the total number of web pages):

$$\text{PageRank}(p) = \frac{1 - d}{N} + d \left(\sum_{i=1}^k \frac{\text{PageRank}(q_i)}{\text{out degree of } q_i} \right)$$

Thus, the PageRank of a page grows with the importance (=PageRanks) of the pages which point to it. An endorsement (=link) from a prominent (high ranking) site, like Yahoo!,¹² contributes to a page's PageRank much more than

¹¹available online at <http://www10.org/cdrom/papers/289/> as of October 2001.

¹²Yahoo! Inc. <http://www.yahoo.com/>

an incoming link from some obscure personal homepage. Our image ranking schemes will imitate this property. In particular, the rankings of images will grow with the importance of the pages which contain them.

2.2 The Mutual Reinforcement Approach

Kleinberg's Mutual Reinforcement approach, introduced in Kleinberg [1999], aims to find hubs and authorities which pertain to a given topic t . The key observation behind the approach is that hubs and authorities which pertain to t display a *mutually reinforcing* relationship: For a page to be considered a good t -hub, it must point to many t -authorities, while a page is considered to be an authority on topic t only if many hubs deem it as such (and point to it). Because of this relationship, prominent t -hubs and t -authorities tend to form *communities*, which can be seen as densely inter-connected bipartite portions of the Web-graph.

The algorithm starts by assembling a collection \mathcal{C} of Web pages, which should contain many high quality Web pages which pertain to a given topic t . It then analyzes the link structure induced by that collection, in order to find the authoritative pages (and the hubs) on topic t .

Denote by q a term-based search query which describes the topic of interest t . The collection \mathcal{C} is assembled as follows:

- A *root set* S of pages is obtained by applying a term based search engine, such as AltaVista, to the query q . This is the only step in which the lexical content of the Web sites is examined.
- From S , a *base set* \mathcal{C} is derived, that consists of (a) pages in the root set S , (b) pages that point to a page in S and (c) pages that are pointed to by a page in S .

The collection \mathcal{C} and its link structure induce a $|\mathcal{C}| \times |\mathcal{C}|$ adjacency matrix, which is denoted by W .

Each page $s \in \mathcal{C}$ is then assigned a pair of weights, a hub-weight $h(s)$ and an authority weight $a(s)$, based on the following two principles:

- The quality of a hub is determined by the quality of the authorities it points at.
- A page is only as authoritative as the quality of the hubs which deem it as such.

The top ranking pages, according to both types of weights, form the Mutually Reinforcing communities of hubs and authorities.

Kleinberg uses the following iterative algorithm to assign the weights:

- (1) Initialize $a(s) \leftarrow 1$, $h(s) \leftarrow 1$ for all pages $s \in \mathcal{C}$.
- (2) Repeat the following operations until convergence:
 - Update the authority weight of each page s (the \mathcal{I} operation):

$$a(s) \leftarrow \sum_{\{x|x \text{ points to } s\}} h(x)$$
 - Update the hub weight of each page s (the \mathcal{O} operation):

$$h(s) \leftarrow \sum_{\{x|s \text{ points to } x\}} a(x)$$
 - Normalize both sets of hub and authority weights.

Note that applying the \mathcal{I} operation is equivalent to assigning authority weights according to the result of multiplying the vector of all hub weights by the matrix W^T . The \mathcal{O} operation is equivalent to assigning hub weights according to the result of multiplying the vector of all authority weights by the matrix W .

Kleinberg showed that this algorithm converges, and that the resulting authority weights [hub weights] are the coordinates of the normalized principal eigenvector¹³ of $W^T W$ [of $W W^T$]. The pages which correspond to the largest coordinates of these eigenvectors are returned by the algorithm as the principal community of authorities[hubs].

The two matrices $W^T W$ and $W W^T$ are well known in the field of bibliometrics:

- (1) $W^T W$ is the *co-citation matrix* [Small 1973] of the collection. $[W^T W]_{i,j}$ is the number of pages which jointly point at (cite) pages i and j .
- (2) $W W^T$ is the *bibliographic coupling matrix* [Kessler 1963] of the collection. $[W W^T]_{i,j}$ is the number of pages jointly referred to (pointed at) by pages i and j .

It is important to note that the outcome of the algorithm, namely the communities of hubs and authorities which the algorithm will identify, is determined solely by the adjacency matrix W of the collection \mathcal{C} . The adjacency matrix implies the co-citation and bibliographic coupling matrices, and it is the eigenvectors of these matrices, in turn, which uniquely determine the principal communities of hubs and authorities.

Our co-citation based image retrieval schemes basically imitate this algorithm, albeit with different adjacency matrices. Defining the adjacency matrices to be used will suffice to uniquely define our schemes.

3. LINK STRUCTURE ANALYSIS FOR FINDING AUTHORITATIVE IMAGES

We now describe our method for finding authoritative images given a query. First, we formally define the image collections which are analyzed. Next, we explain how such collections are assembled from a given query. Finally, we present our image ranking schemes.

3.1 Formal Definition of the Model

A page p is said to contain an image i (denoted by $p \rightsquigarrow i$) in either of the following two cases:

- (1) p displays i : When page p is loaded in a Web browser, i is displayed. For example, either of the following html directives:
`` or
``.
- (2) p points to i 's image file (in some image file format such as .gif or .jpeg). For example,

¹³The eigenvector which corresponds to the eigenvalue of highest magnitude of the matrix.

`nice image`. Note that p does not contain i when p points to an HTML file which contains i , even if i is the only visible object in the HTML file.

We define a topical WWW *image collection* as a quadruple $\mathcal{IC} = (\mathcal{P}, \mathcal{I}, \mathcal{L}, \mathcal{E})$, where \mathcal{P} is a set of Web pages (many of which deal with a certain topic t), \mathcal{I} is the set of images which are contained in \mathcal{P} , $\mathcal{L} \subseteq \mathcal{P} \times \mathcal{P}$ is the set of (directed) links which exist on the Web between the pages of \mathcal{P} , and $\mathcal{E} \subseteq \mathcal{P} \times \mathcal{I}$ is the relation *page p contains image i* .

We denote by W the adjacency matrix of the page-to-page relation \mathcal{L} , and by $M = [m_{ij}]$ the $|\mathcal{P}| \times |\mathcal{I}|$ adjacency matrix of the page-to-image relation \mathcal{E} .

Page-Image Adjacency. Two of the most important observations of link-structure analysis are the following:

- (1) The notion of authority being conferred through links from a pointing resource to a pointed resource. In our context, the pointing resources are Web pages, while the pointed resources are the images.
- (2) The topical similarity between resources (images, in our case) which is inferred through co-citation.

Both principles are reflected in the adjacency relation which exists in the data. The adjacency matrix conveys the flow of authority, while the entries of the co-citation matrix, which the adjacency matrix implies, define the strength of the topical affinities between the resources. We therefore aim to define an adjacency relation between Web pages and images, in a manner which best reflects both the flow of authority from pages to images, and the topical affinities between the various images.

There are several reasonable definitions for such adjacency relations. The intuition behind these definitions is perhaps best explained through an example. Consider the scenario depicted in Figure 1, which consists of five Web pages P_1, \dots, P_5 and four images, one of them replicated.

The most basic adjacency relation which comes to mind is to adopt the page-to-image relation \mathcal{E} , defined above, as the adjacency relation (and M as the adjacency matrix). M represents the outright manner for a page to endorse an image, which is simply to display it or point to it. This approach also reflects some topical affinities between images, through the corresponding co-citation matrix $M^T M$. Pairs of images which are co-contained in the same page are considered topically related. For example, note the entry which corresponds to the runner and tennis player in Figure 2.

This approach, however, fails to convey some other fairly intuitive topical relations, such as between the soccer player and the tennis player. These images appear in pages which are co-cited; assuming we agree that co-cited pages are topically related, then perhaps the images which are contained in them are also. To reflect such a connection, we need to consider the adjacency matrix WM , which associates each page p with the images that are displayed in *pages to which p links*. Using WM as the adjacency matrix, the co-citation matrix

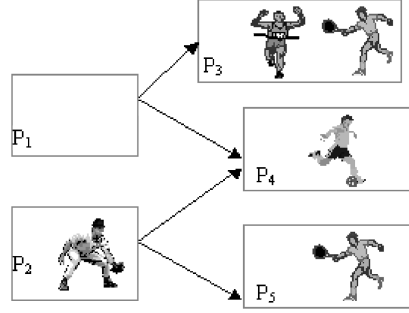


Fig. 1. An example page/image case.

P ₁	0	0	0	0
P ₂	0	0	0	1
P ₃	1	0	1	0
P ₄	0	1	0	0
P ₅	1	0	0	0

(a) M

P ₁	2	0	1	0
P ₂	0	1	0	0
P ₃	1	0	1	0
P ₄	0	0	0	1

(b) $M^T M$

Fig. 2. The $k = 0$ adjacency and co-citation matrices for the example case.

P ₁	1	1	1	0
P ₂	1	1	0	0
P ₃	0	0	0	0
P ₄	0	0	0	0
P ₅	0	0	0	0

(a) WM

P ₁	2	2	1	0
P ₂	2	2	1	0
P ₃	1	1	1	0
P ₄	0	0	0	0

(b) $(WM)^T WM$

Fig. 3. The $k = 1$ adjacency and co-citation matrices for the example case.

$M^T W^T WM$ (Figure 3) now reflects some topical affinity between the soccer image and the tennis image.

This approach also suffers from some obvious setbacks. The affinity between the soccer image and the runner image is considered as strong as the affinity between the tennis image and the runner, although it seems logical that the tennis and runner images are more tightly coupled, since they appear in the same page. A greater problem exists with the connection between the image of the baseball player and the images of the soccer and tennis players. Why is a linkage between the soccer image and the tennis image inferred by P_2 co-citing pages P_4 and P_5 , while no linkage is inferred between those two images and the baseball image, contained in P_2 itself? Perhaps the answer to these issues lies in using the matrix $(W + I_{|\mathcal{P}|})M$ as the adjacency matrix (where $I_{|\mathcal{P}|}$ is the $|\mathcal{P}| \times |\mathcal{P}|$ identity matrix).

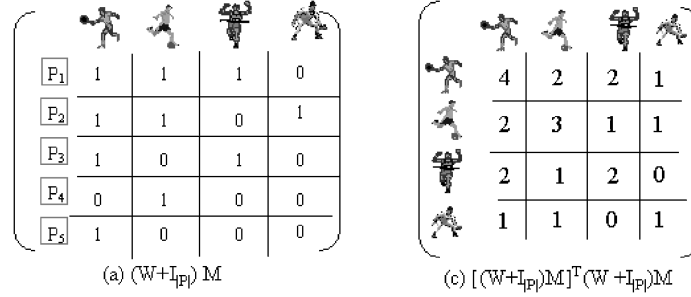


Fig. 4. Two times the $k = 0.5$ adjacency matrix (and the corresponding co-citation matrix) for the example case.

The matrix $(W + I_{|P|})M$ defines each page to be adjacent both to the images which are contained in it, and to the images which are contained in pages to which it points. The co-citation measure which this adjacency matrix implies (the corresponding co-citation matrix is $[(W + I_{|P|})M]^T(W + I_{|P|})M$), now addresses the concerns raised in the previous case (as shown in Figure 4).

A closer look at our three proposed adjacency matrices (M , WM and $(W + I_{|P|})M$) reveals that they are all members of the following parametric family of adjacency matrices (up to, perhaps, a constant factor):

$$A_{IC} = \{[kW + (1 - k)I_{|P|}]M : 0 \leq k \leq 1\}$$

By denoting $A_{IC}(k) \triangleq [kW + (1 - k)I_{|P|}]M$, we have $M = A_{IC}(0)$, $WM = A_{IC}(1)$, and $(W + I_{|P|})M = 2 \cdot A_{IC}(\frac{1}{2})$. In general, choosing large values of k will introduce bias towards relations between pages and images contained in pages linked from them, while small values of k will boost the relationship between pages and the images that they themselves contain.

Our experiments and sample results were derived using the three adjacency matrices defined above, although we do not claim that any of these choices is in any way optimal.

3.1.1 Weighted Relations. The definitions of the previous subsection are easily extended to the case where both \mathcal{L} and \mathcal{E} are weighted relations, that is $\mathcal{L} \subseteq \mathcal{P} \times \mathcal{P} \times \mathbf{R}^+$,¹⁴ is the set of weighted page to page links and $\mathcal{E} \subseteq \mathcal{P} \times \mathcal{I} \times \mathbf{R}^+$ is the weighted page-image relation. Weighted relations are derived by assigning weights to the links (relations) which reflect the amount of authority that the pointing (containing) page confers to the pointed page (image). Possible factors which may contribute to the weight of a link or relation include the following (the first factor is considered in PicASHOW):

—Anchor text which is relevant to the query. Such text around a link raises our confidence that the pointed page or contained image is relevant to the topic at hand [Chakrabarti et al. 1998]. Similarly, in the case of the page-image relation, when the content of the ALT field is relevant to the query, then the image is most likely related to the topic of interest.

¹⁴ \mathbf{R}^+ is the set of non-negative real numbers.

- The position of the link[image] in the pointing [containing] page. Many search engines consider the text at the top of a page as more reflective of its contents than text further down the page. The same line of thought can be applied to the links/images which appear in a page, with those which are closer to the top of the page receiving more weight than those appearing at the bottom of the page.

3.2 Assembling a Topical Collection

Our assumption in applying link-structure analysis when searching for quality images on topic t is that t -relevant pages will contain quality images on t . Thus, by examining a large enough set of t -relevant pages, we should be able to identify high quality t -images. Therefore, the first step in assembling a topical collection of images is to assemble a large collection of t -relevant pages. This collection is assembled in the same manner as described in Section 2.2. That is, for a query q which describes the topic t , we assemble a q -induced collection of Web pages by submitting q first to traditional search engines, and adding pages that point to or are pointed by pages in the resultant set. This provides us with the page set \mathcal{P} and the page-to-page link set \mathcal{L} . Note that we do not utilize any image search engine in this step—we are using only standard (HTML-page finding) search engines. Once we compile the page set \mathcal{P} , we define the set \mathcal{I} as the set of images which are contained in \mathcal{P} (this also implies the page-to-image relation \mathcal{E}).

This scheme for assembling the topical image collection can be implemented with reasonable overhead, in terms of both computations and storage, by many standard WWW search engines. All search engines continuously crawl the Web with robots, which collect the textual content of the pages to be indexed. Many engines (such as AltaVista, Google and Lycos¹⁵), in addition, collect connectivity information that captures the information regarding the links between the pages as they crawl the Web. Our scheme requires the engines to also catalog, for each page, which images are contained (or pointed to) in the page. However, the images themselves need not be stored. As explained below, each image only requires the storage of a 32 byte signature, plus a URL where it can be found.

When building the image collection $\mathcal{IC} = (\mathcal{P}, \mathcal{I}, \mathcal{L}, \mathcal{E})$, we must consider a common authoring technique of Web pages. Specifically, when a Web site creator encounters an image of his liking on a remote server, the usual course of action would be to copy the image file to the local server, thus replicating the image. The motivation behind this practice is to enable the author's page to load faster from within the author's domain/organization, since the displayed images are stored locally.

This behavior of authors with respect to images is different from their corresponding behavior with respect to HTML pages. In most cases, authors will not copy a remote page (or some portion of its contents) to the local servers; rather, they provide links from their site to the remote page. There are exceptions to this rule (such as the replication of system manuals or software APIs), but for most types of content, HTML pages are not replicated. This authoring mode

¹⁵Lycos Inc. <http://www.lycos.com/>

has two important implications for link-based image search techniques, which are in contrast to the corresponding link based techniques for searching Web pages. We expand on these now.

Identifying replicated images. We must identify that multiple pages contain a certain image, even when the pages contain different copies of the image. Thus, images cannot be identified by their URIs, but must be identified by their content. In contrast, when applying link-analysis in search of authoritative pages, identifying replications is less crucial. Satisfactory results can be obtained even when the issue of page-replications is ignored.

Fortunately, it is possible to decide whether two images are replicated, with a relatively high probability, by examining a small portion of the image. In PicASHOW, we only download the first 1024 bytes of the image and apply a double hash function to these bytes, so that each image is represented by a signature consisting of 32 bytes. Two images with the same signature are considered identical. Our experience shows that very rarely do different images result in the same signature since the first 1024 bytes usually capture the header information as well as the first part of the image itself. The storage overhead which is associated with each image is thus quite minimal. Note that replications of the same image result in only one 32-byte signature (and one URL) in terms of storage requirements.

Filtering non-informative images. Link-analysis based page-search methods usually interpret a link from page p to page q as a measure of authority which p confers on q [Kleinberg 1999]. However, there are many kinds of links which confer little or no authority [Chakrabarti et al. 1999]; we refer to these as *non-informative* links. Some examples for such links are intra-domain (inner) links (whose purpose is to provide navigational aid in complex Web sites), commercial/sponsor links, and links which result from link-exchange agreements. A crucial task which should be completed prior to analyzing the link structure of a given collection, is to filter out as many non-informative links as possible.

Similarly, filtering non-informative page-to-image links is crucial for successful link-based image retrieval. Site banners and logos can be thought of as the image equivalents of non-informative links. These images introduce a large amount of noise into image collections, which we would like to be able to filter. When building the set of page-to-page links \mathcal{L} , we identify and filter out intra-domain links, ruling them to be navigational links which do not confer authority. However, the practice (described above) of image replication implies that filtering out the intra-domain page-to-image links of \mathcal{E} will be destructive as we may also lose quality images in this fashion. We thus introduce a few heuristics, that can mitigate the noise that is introduced by non-informative images:

- Banners and logos tend to be wide and short. We can thus filter out images with an aspect ratio greater than some threshold. Note that we only need to examine the image header for this information.
- Images which are stored in small files (less than 10 kilobytes, for example) tend to be banners. Even if they are not banners, they are usually not quality topical images. Therefore, they can be filtered from the collection.

- Images that are stored in files whose names contain the words *logo* or *banner* are probably logos and banners.
- In addition to banners and logos, people tend to include other non-informative images such as clipart in order to liven up their Web page. These include colorful buttons, bars, mail boxes, spinning globes, and so forth. Some of these are highly popular, and are replicated and used in large numbers on the Web. We consider these the equivalent of stop words in information retrieval [van Rijsbergen 1979] and thus term them *stop images*. Many of these stop images are filtered based on the aspect ratio and file size heuristics. In addition, we have constructed a list of stop images (common names of some of these images, and 32-byte signatures of others), and we filter out any image that appears in this list. Currently, this list is assembled manually. It is, however, feasible to compile such a list based on distribution statistics of images on the WWW.

These heuristics do not filter out all the noise caused by non-informative images. Some non-informative images survive this process, and introduce noise into our page-to-image adjacency matrices. We found that this noise affects the image rankings more than usually happens in the corresponding page ranking schemes (where page-to-page adjacency matrices are used). As a consequence, link-based image search seems to be more noisy than link-based page search, and specifically may be easier to spam. Devising more elaborate and effective filtering schemes is left for future work.

3.3 Image Ranking Schemes

After assembling an image collection $\mathcal{IC} = (\mathcal{P}, \mathcal{I}, \mathcal{L}, \mathcal{E})$ pertaining to a certain topic t , we need to rank the images of \mathcal{I} with respect to t . We assume that every page $p \in \mathcal{P}$ is associated with a t -relevance score, denoted by $r_t(p)$. Note that this is not a limiting assumption, since we can always calculate the authority scores of the pages and use them as relevance scores. In particular, the collection \mathcal{IC} contains the linkage information which is required to calculate authority scores by the mutual reinforcement approach.

Below is a list of the ranking schemes with which we have experimented. They are divided into three categories: A naive image in-degree approach, PageRank-influenced ranking schemes, and co-citation based analyses. These ranking schemes are based on the matrices which were defined in Section 3.1.

- (1) *In-degree rank according to the matrix M* . Here, the score of image i equals $\sum_{\{p \in \mathcal{P} \mid p \text{ contains } i\}} m_{p,i}$, where $m_{p,i}$ is the weight associated with the page-image relation $p \rightsquigarrow i$. That is, the score of image i is the sum of the weights of all relations $p \rightsquigarrow i$ for all pages p which contain image i in the collection.
- (2) *PageRank [Brin and Page 1998] influenced ranking schemes*. Under the hypothesis that images which are contained in t -relevant pages should be of higher quality (with respect to t) than images contained in t -irrelevant pages, we factor the relevance score of page p ($r_t(p)$) into the score of image

- i . In particular, we set the score of image i to equal $\sum_{\{p \in \mathcal{P} \mid p \text{ contains } i\}} r_t(p)$. In the case that M is a weighted matrix (and not simply a binary matrix), the straightforward variant of this score is $\sum_{\{p \in \mathcal{P} \mid p \text{ contains } i\}} r_t(p) m_{p,i}$.
- (3) *Co-citation based analyses*. Each of the matrices M , WM and $(W + I_{|\mathcal{P}|})M$ may serve as the citation (adjacency) matrix. The corresponding co-citation matrix $M^T M$, $(WM)^T(WM)$ or $[(W + I_{|\mathcal{P}|})M]^T[(W + I_{|\mathcal{P}|})M]$ is used for the purpose of the analysis. Specifically, we tested the rankings that are produced by the mutual reinforcement approach and by SALSA.¹⁶ We can also use the t -relevance score $\{r_t(p), p \in \mathcal{P}\}$ to enhance our co-citation analysis and boost the rankings of images that are cited by highly t -relevant pages. As an example, consider applying co-citation analysis to the $|\mathcal{P}| \times |\mathcal{P}|$ page-to-image matrix M_R that is defined as follows:

$$[M_R]_{i,j} \triangleq m_{i,j} \sqrt{r_t(i)}$$

By examining the t -relevance image co-citation matrix $M_R^T M_R$, we note that when M is unweighted (a binary adjacency matrix), $[M_R^T M_R]_{i,j}$ sums the relevance weight of all pages which co-display images i and j :

$$[M_R^T M_R]_{i,j} = \sum_{\{k: k \sim i, k \sim j\}} [M_R]_{k,i} [M_R]_{k,j} \quad (1)$$

$$= \sum_{\{k: k \sim i, k \sim j\}} (\sqrt{r_t(k)})^2 = \sum_{\{k: k \sim i, k \sim j\}} r_t(k) \quad (2)$$

Thus, not all image co-citations are considered equal—highly relevant pages endorse their co-contained images to a larger extent than do less relevant pages.

4. IMAGE CONTAINERS AND IMAGE HUBS

One of the major benefits of hyperlink-based image search is that in addition to finding good images which pertain to a certain query, it can identify Web pages that are rich in relevant images, or from which many images are readily accessible. Our proposed co-citation based ranking schemes naturally allow for such Web pages to be found.

While we have concentrated, in the previous section, on how to rank the authoritative images in \mathcal{IC} , we can similarly find Web pages whose role corresponds to that of hubs. Just as hubs were defined as pages which link to many authoritative pages, in our context *image hubs* should be pages which are, in some sense, linked to many authoritative images.

However, the notion of an *image hub* is somewhat ambiguous. Do we, by calling p an “image hub”, mean that many authoritative images are displayed in p , or do we mean that p points to many pages which contain quality images? We claim that both possible interpretations are of value, and so we

¹⁶An in-depth description of SALSA [Lempel and Moran 2000] is beyond the scope of this paper. However, SALSA is based on co-citation, and its analysis, just like Kleinberg’s mutual reinforcement approach, is completely governed by the adjacency matrix that is used.

define them separately as follows: Pages which contain high-quality images are called *image containers*, while pages which point to good image containers are called *image hubs*. Thus, image hubs are once removed from the authoritative images themselves, which are contained (as the name implies) in the image containers.

Our co-citation based image retrieval schemes can find both image containers and image hubs, either separately or in some mixed manner. The outcome depends on the type of adjacency matrix used to describe the collection \mathcal{IC} , which, in turn, implies the bibliographic coupling matrix which governs the ranks of pages as image hubs/containers (The technical details of how hub ranks are derived from the adjacency matrix were given in Section 2.2).

When using the adjacency matrix M (or M_R), the p 'th row describes which images are contained in page p . The pages whose coordinates will stand out in the principal eigenvector of the matrix MM^T will, accordingly, form a community of image containers. Image hubs are likely to be found when using the adjacency matrix WM , whose p 'th row describes which images are contained in pages to which p links (the corresponding bibliographic coupling matrix will be WMM^TW^T). Co-citation analysis using the matrix $(W + I_{|P|})M$ allows us to find communities of pages which both contain images and which point to other image containers, since the p 'th row there details the images which are contained either in p or in pages to which p links. In general, when using the adjacency matrix $A_{\mathcal{IC}}(k)$, high values of k shift the analysis towards image hubs, while lower values of k accentuate image containers.

One of the most striking image containers found throughout our experiments with PicASHOW is the fractal image container <http://sprott.physics.wisc.edu/fractals.htm>. This page contains over a hundred images of fractals. Another example is the image container for the query *Magritte*, <http://www.xs4all.nl/~renebos/magritte.html>. The page itself contains just 5 words: Rene Magritte : 10 Genius Paintings. Underneath that phrase are displayed ten of Magritte's masterpieces.

Here are a couple of examples for image hubs:

- (1) For the query *pyramids*, many relevant images are just a click away from the following URLs: (a) The Ancient Egypt page, pyramids section at <http://members.aol.com/TeacherNet/AncientEgypt.html> and (b) The Art History Resources: Part 2 Ancient Art page, Egypt section at <http://witcombe.sbc.edu/ARTHancient.html#AncEgypt>
- (2) For the query *Yosemite*, "Photos by Rick Ellis—Yosemite" (<http://www.fnet.net/~ellis/photo/yosemite.html>) is both an image container and an image hub. Actually, it is a Yosemite hub in general, since it points to many valuable resources on Yosemite National Park. Many Yosemite images are accessible from http://home.earthlink.net/~mgordon324/sierra_new.htm, Michael Gordon's High Sierra page, which is an image container and hub for the Sierra Nevada mountain range (which includes the Yosemite National Park).

Table I. Rank Method/Category Precision@10 Values

Category	ID	WPR	M_R	WM	$(W + I)M$	S- M	S- $(W + I)M$
Art	5.67	8	8	10	10	7.33	9
Nature	6	8	8	8.33	8.33	8	8
Celebrities	7	7	7	7.33	7	8	7
General locations	5	7.5	8	6.5	6	6	6.5
Transportation	4	7	7.5	7	7	5.5	6
Children	6	7	6	5.33	6	6	6
Specific locations	5	6	6	4	4	3.5	7
Concepts	1.33	3	4.33	4.67	4.33	3.33	2.67

5. EVALUATION OF RESULTS

In the course of our work on PicASHOW we obtained results for many queries. We tested the performance of several of our ranking schemes on 25 queries from 8 diverse categories. Table I reports the average precision@10 values which the tested ranking schemes attained.

Throughout these experiments, the relations \mathcal{L} and \mathcal{E} were weighted (extra weight was given to links with anchor text [ALT field text] which resembled the query). Following are the ranking schemes that were tested (recall the definitions of the ranking schemes from Section 3.3):

- (1) ID: In-degree rank of images, according to the matrix M .
- (2) WPR: The (weighted) PageRank influenced scheme, using the weighted matrix M and the relevance of the containing pages to the query t .
- (3) M_R : Co-citation analysis with the mutual reinforcement approach, using the adjacency matrix M_R .
- (4) WM : Co-citation analysis with the mutual reinforcement approach, using the adjacency matrix WM .
- (5) $(W + I)M$: Co-citation analysis with the mutual reinforcement approach, using the adjacency matrix $(W + I_{|P|})M$.
- (6) S- M : Co-citation analysis with SALSA, using the adjacency matrix M .
- (7) S- $(W + I)M$: Co-citation analysis with SALSA, using the adjacency matrix $(W + I_{|P|})M$.

Here are the queries which comprised each category:

- Art: “Vincent Van Gogh”, “Magritte”, “Roman architecture”.
- Nature: “fractals”, “solar system”, “snakes”.
- Celebrities: “Bill Gates”, “Meg Ryan”, “Michael Jordan + Bulls”.
- General locations: “Paris”, “Crete”, “Yosemite”, “Kilimanjaro”.
- Transportation: “nuclear submarine”, “Volkswagen beetle”, “race car”, “Concorde aircraft”.
- Children: “Pickachu”, “moomin”, “Cheshire cat”.
- Specific Locations: “Eiffel Tower”, “Golden Gate Bridge”.
- Concepts: “anthropology”, “jogging”, “early aircrafts”.

Additional findings which complement Table I:

- The in-degree ranking scheme is almost constantly out-performed by the more complex ranking schemes.
- Most ranking schemes did not cope well with finding images of specific locations. Part of this failure was due to *topic drift*, or *diffusion* [Kleinberg 1999; Bharat and Henzinger 1998]. For example, many of the returned images for the “Golden Gate Bridge” query were either general *San Francisco* images, or pictures of famous bridges.
- Our methods failed to find authoritative images for the “concepts” queries. With “early aircrafts”, most methods could not differentiate between aircrafts from different eras; “jogging” results were a mixed array of images of outdoor activities; and the “anthropology” results included buildings of anthropology faculties.
- The results returned by schemes WPR and M_R were often quite similar, and the same is true for the results of schemes WM and $(W + I)M$. For the first pair, this is probably due to the strong influence of the page relevance scores $r_t(p)$. Authority scores are usually very skewed, with most pages having very low scores. In the first pair of ranking schemes, the images in non-authoritative pages are practically ignored, and only the images of the relatively few authoritative pages are considered. As for the second pair, our assumption was that using the adjacency matrix $(W + I_{|P|})M$ would produce results that were a mixture of the results attained with the matrices M and WM . It turns out that the results of WM are much more dominant, hinting that in order to get hybrid results one should try using adjacency matrices of the form $(W + \lambda I_{|P|})M$ with $\lambda \gg 1$.
- For 20 of the 25 queries tested, at least one ranking scheme had at least 7 good results in its top-10 images. If we ignore the “concepts” category, and examine the two best ranking schemes for every remaining query, we find that all 22 queries returned at least 12 relevant images (out of a possible 20). This means that a collage of the (distinct) images in the returned top-10 lists will usually satisfy user needs. Such a collage, for “difficult” queries, might contain many irrelevant images. However, achieving high precision is less crucial in image search applications than it is for page searches. Users who are presented with thumbnails of 50 images can easily filter out the irrelevant images and concentrate on the relevant ones. The cost of imprecision is far greater for page searches, where digesting returned lists of URLs is much more cumbersome.

To further demonstrate the abilities of PicASHOW, images of sample results were shown in Lempel and Soffer [2001], sometimes alongside results of commercial image search engines on the same query. We cannot show those images here; however, we show the URLs of those images (see Appendix A), provide some details on the ranking scheme that was applied in each case, and highlight other noteworthy points that follow from each example.

- (1) Table II displays the URLs of PicASHOW's results for the query "*Michael Jordan*". The image collection consisted of 1083 images, and the results were derived by applying the PageRank-influenced ranking scheme. The URLs of the retrieved images exemplify the practice of image replication among Web authors. Note that two of those images are animated GIFs.
- (2) Table III has the URLs of our results for the query "*Jaguar car*". The image collection in this case was very small, and consisted of just 67 images. The ranking scheme was Kleinberg's algorithm, using the adjacency matrix M . Note that while PicASHOW's results all come from different servers, most of the Lycos images are from a single server.
- (3) Table IV displays the URLs of PicASHOW's results for the query "*Kilimanjaro*". The 309 images of this example were ranked by SALSA, using the adjacency matrix M . Note how some of the image names do not contain anything resembling the query, but rather the name "*Kibo*", which is the name of one of Mt. Kilimanjaro's peaks.
- (4) Table V displays the URLs of PicASHOW's results for the query "*Vincent Van Gogh*". The image collection consisted of 582 images, and the results were derived by applying Kleinberg's algorithm using the adjacency matrix $(W + I_{|\mathcal{P}|})M$.
- (5) Table VI displays the URLs of our results for the query "*Solar System*". The image collection consisted of 682 images, and the images were ranked by weighted in-degree according to the relation \mathcal{E} .

6. DISCUSSION AND FUTURE WORK

Our experiments with PicASHOW and the examples showcased throughout [Lempel and Soffer 2001] demonstrate that PicASHOW, while relying on very little besides link analysis, demonstrates retrieval abilities comparable to those of available WWW image search engines. In addition, PicASHOW's retrieval of image containers and hubs is a natural and useful extension of the image search paradigm, which, to the best of our knowledge, has not been previously pursued.

As in the case of link analysis in Web page search, PicASHOW performs best on *wide-topic* queries, where the search topic is of wide interest on the Web. Queries on obscure topics of little interest will most likely fail to produce quality results. Results of queries on specific aspects of wide topics will often *drift*, returning images that are relevant to the general topic, not necessarily to the requested aspect of it. In addition, our ranking techniques are most effective on topics where image replication is likely to occur. For example, people are likely to display replications of famous works of art and of publicly released images of celebrities. Natural landmarks may also induce image replications. However, a query such as *Paris* is less likely to achieve quality results, since although the topic of "Paris" has broad presence on the Web, people will often display images of their own vacation in Paris, rather than replicated versions of the city's landmarks.

Since PicASHOW performs no image analysis whatsoever, it cannot handle queries that contain image qualifiers such as color, orientation, and other

specific features. For example, PicASHOW can retrieve images of Michael Jordan, but not of Michael Jordan *wearing a suit*. It can find images of Jaguar cars, but not of *red* Jaguar cars, and while it will nicely rank images of Mount Kilimanjaro, one cannot ask for those images *not to contain trekkers*, or to be taken *from below*. Note however, that link-analysis based techniques could still be of value for such queries. For example, PicASHOW could be used as an initial filter to find candidate images on the topic of interest (e.g., Jaguar cars). Some form of image analysis could then be performed on these candidate images in order to select those that further satisfy the image qualifications (e.g. which of the resulting Jaguar cars are red).

It is obvious from the above discussion that link analysis by itself is not a silver bullet for Web image retrieval, and should be augmented with other retrieval means. Setting aside content-based algorithms, at least two traditional IR techniques should be integrated into a Web image retrieval system:

- Text-based retrieval, with an emphasis on sophisticated methods of automatically associating descriptive text with each image.
- An effective (query-specific) page-ranking algorithm (Since the relevance of the displaying page with respect to a query is a good indicator of the quality images which it displays).

Several interesting extensions of the image search schemes are feasible:

- (1) In Section 2 we described the mutual reinforcement approach [Kleinberg 1999], which ranks pages with respect to specific queries, and PageRank [Brin and Page 1998], which assigns each Web page a global importance measure. PicASHOW follows the former, and ranks images in topical, query-specific image collections which it assembles. It is challenging to suggest, implement and evaluate a global image ranking scheme, analogous to PageRank.
- (2) Kleinberg [1999] demonstrated that non-principal communities of authoritative pages can distinguish between topics in multi-topic collections, and between pages which present opposing views on polarized topics (such as the *pro-life* and *pro-choice* views on *abortion*). It thus seems interesting to investigate the non-principal communities of images which arise from the various proposed co-citation measures.
- (3) Most CBIR systems support image similarity queries, and extending our approach to support such queries would enhance its appeal. Many algorithms have been proposed for finding pages related to a sample page (or set of pages) on the Web by link analysis [Dean and Henzinger 1999; Aridor et al. 2000]. The main idea is to grow a Web-graph around the given seed pages, and then find the dominant authorities in that graph. It seems possible to adapt PicASHOW in the same fashion to support similar-image queries. The input to such queries can be either URLs of sample images, or URLs of sample image containers.
- (4) In our current prototype, we define the images in an image collection $\mathcal{IC} = (\mathcal{P}, \mathcal{I}, \mathcal{L}, \mathcal{E})$ to be the images contained in the set of pages \mathcal{P} . Recall

that \mathcal{P} was actually a neighborhood of pages around a set of root pages S on some topic t . Consider the set of root images \mathcal{I}_S , defined as the set of images that are contained in the pages of S . We currently expand \mathcal{I}_S into the final set of images \mathcal{I} by following page-to-page links, in other words by expansion on the *page plane*. However, \mathcal{I}_S may also be expanded by adding images that are contained in other pages containing replications of the root images. The premise is that, if pages that are not linked to the root set of pages S , contain replications of images from \mathcal{I}_S , they may also contain other images of relevance to t . We term as *image plane expansion* the process of adding those pages, and their contained images, to \mathcal{I} . Technically, such expansion requires information of the form “which pages contain the following (signature of an) image”. This information is not currently available on the Web, but search engines which will collect page-to-image connectivity information can support such queries, and thus enable image plane expansion as well.

APPENDIX

A. URLS OF IMAGES

For the URLs of PicASHOW's results, multiple URLs are given for replications of the same image, when applicable.

Table II. URLs of *Michael Jordan* Images

URLs PicASHOW's " <i>Michael Jordan</i> " Images
http://views.vcu.edu/abaididi/jordan01.gif http://www.geocities.com/Colosseum/Sideline/1534/jordan01.gif http://acnc.spa.k12.mi.us/powers/jordan01.gif http://www.eng.fsu.edu/toliver/jordanmovie.gif http://www.2.gvsu.edu/%7Ejirtlee/Bettermovingdunk.gif http://aesd.sk.ca/scp/images/AIR_JORDAN.gif (animated gif)
http://www.geocities.com/Colosseum/Sideline/1534/jumper.jpg http://scnc.sps.k12.mi.us/powers/jumper.jpg http://icdweb.cc.purdue.edu/fultona/MJ11.jpg
http://www.engin.utnd.umich.edu/jafreema/mj/mjpgs/jordan5-e.gif http://www.angelfire.com/ny/Aaronakickasspages/images/1-11.JPG http://homepages.cu.rmit.edu.au/dskiba/mjsmile.jpg http://acnc.sps.k12.mi.us/woolworl/jordan.jpg http://www.fidelweb.com/graphic/jordan4.jpg
http://www.engin.umd.umich.edu/jafreema/mj/mjpgs/jordan10-e.jpg http://www.engin.umd.umich.edu/jafreema/pictures/1991.gif http://www.metal.chungnarn.ac.kr/myoungho/1991.gif (animated gif)
URLs of Scour's " <i>Michael Jordan</i> " images
http://www.geocities.com/SunsetStrip/2546/mjjuwan.jpg http://www.geocities.com/SunsetStrip/2546/jordanmvp2.jpg http://www.unc.edu/lbrooks2/jordan2.jpg http://www.geocities.com/Colosseum/Tarck/7823/JORDAN_ALLSTAR_1.JPG http://www.big.du.se/joke/f1-96/pics/car/jordan96_car.jpg http://www.unc.edu/lbrooks2/mjbugs.jpg

Table III. URLs of *Jaguar Car* Images

URLs of PicASHOW's <i>Jaguar Car</i> Images
http://www.classicar.com/museums/weishjag/outside.gif
http://www.ferrari-transmissions.co.uk/home2.jpg ¹⁷
http://www.jtc-nj-com/Doylestowncrowd.jpg
http://www.jaguar-association.de/images/verkaufabilder/12-00-teach/rs100s-lg.jpg
http://www.j-c-c.org.uk/images/drive.jpg
http://www.seattlejagclub.org/IMAGES/picyak.jpg
URLs of the Lycos <i>Jaguar Car</i> Images
http://www.auto.com/art/reviews/98_jaguar_xjr/98_jaguar_XJR-Interior.jpg
http://highway-one.com/Images/Photos/Jaguar/LaGrassaJaguar4.jpg
http://highway-one.com/Images/Photos/Jaguar/LaGrassaJaguar2.jpg
http://highway-one.com/Images/Photos/Jaguar/LaGrassaJaguar.jpg
http://highway-one.com/Images/Photos/Jaguar/LaGrassaJaguar3.jpg from

Table IV. URLs of *Kilimanjaro* Images

URLs of PicASHOW's " <i>Kilimanjaro</i> " Images
http://www.calle.com/carl/brett.kili.jpg
http://www.premier.org.uk/graphics/programmes/kili001.jpg
http://www.sfusd.edu/cj/kibo.jpg
http://www.nisua.sfusd.k12.ca.us/cj/kibo.jpg
http://www.geocities.com/Yosemite/1015/kili1.jpg
http://seclab.ca.ucdavis.edu/wee/images/kili-summit.gif
http://www.geocities.com/Yosemite/1015/kili2.jpg
http://www.picton-castle.com/jpg/Kilimanjaro_masai.T.jpg
http://www.adventure.co.se/ISTPAGEOFKIBO.jpg

Table V. URLs of *Vincent Van Gogh* Images

URLs of PicASHOW's " <i>Van Gogh</i> " Images
http://www.vangoghgallery.com/images/small/0612.jpg
http://www.scf.uae.edu./wrivera/vangogh.jpg
http://www.openface.ca/vangogh/images/small/0627.jpg
http://www.vangoghgallery.com/images/small/0627.jpg
http://www.sd104.a-cook.k12.il.us/rhauser/vangoghsel.jpg
http://www.openface.ca/vangogh/images/small/0627.jpg
http://www.vangoghgallery.com/images/intro/1530.jpg
http://www.openface.ca/vangogh/images/intro/1530.jpg
http://www.bc.edu/bc-org/avp/cas/fnart/art/19th/vangogh/vangoghself3.jpg
http://sunsite.unc.edu/wm/paint/auth/gogh/entrance.jpg
http://www.ibiblio.org/wm/paint/auth/gogh/entrance.jpg
http://www.southern.com/wm/paint/auth/gogh/entrance.jpg
http://www.bc.edu/bu_org/avp/cas/fnart/art/19th/vangogh/vangogh_starry1.jpg
URLs of Alta Vista's " <i>Vincent Van Gogh</i> " Images
http://www.ElectronicPostcards.com/pc/pics/van12b.jpg
http://www.ElectronicPostcards.com/pc/pics/van5b.jpg
http://www.ElectronicPostcards.com/pc/pics/van1b.jpg
http://www.culturekiosque.com/images5/van.jpg
http://www.ElectronicPostcards.com/pc/pics/van6b.jpg
http://www.ElectronicPostcards.com/pc/pica/van2b.jpg

¹⁷When enlarged, this image reads "Michael Ferrari is my name, but Jaguars are my game". Mr. Ferrari claims to be an independent Jaguar transmission specialist.

Table VI. URLs of *Solar System* Images

URLs of PicASHOW's " <i>Solar System</i> " Images
http://oposite.atsci.edu/pubinfo/jpeg/M16Full.jpg http://www.geoaci.unc.edu/classes/Geo120/SNsmall.gif ¹⁸ http://www.geoaci.unc.edu/classes/Geo15/SNsmall.gif http://nssdc.gafc.nasa.gov/image/planetary/solar_system/family_portraits.jpg
http://www.aeds.org/nineplanets/nineplanets/gif/SmallWorlds.gif http://www.4.net.netage.com/chwu/images/solar_system/nineplanets/SmallWorlds.gif http://www.physics.louisville.edu/tnp/gif/SmallWorlds.gif http://img.iln.net/image/main/astronomy/gif/SmallWorlds.gif http://www.hpcc.astro.washington.edu/mirrors/nineplanets/gif/SmallWorlds.gif http://seds.lpl.arizona.edu/nineplanets/nineplanets/gif/SmallWords.gif
http://www.seds.org/nineplanets/nineplanets/Nineplanets.jpg http://kiss.uni-ij.ai/k4fg0152/devetplanetov/xslike/9planetov.x.jpg http://www.physics.louisville.edu/tnp/NinePlanets.jpg http://img.iln.net/images/main/astronomy/NinePlanets.jpg http://www.hpcc.astro.washington.edu/mirrors/nineplanets/NinePlanets.jpg http://seds.lpl.arizona.edu.edu/nineplnets/nineplanets/NinePlanets.jpg
http://www.solarviews.com/images/rocketvision.gif (animated gif) http://nssdc.gafc.nasa.gov/image/planetary/solar_system/solar_family.jpg
URLs of Ditto's " <i>Solar System</i> " Images
http://www.festivale.webcentral.com.au/shopping/art.com/SYST.jpg http://www.coseti.org/images/12358.jpg http://www.greenbuilder.com/sourcebook/SourcebookGifs/HeatCoolSolar.2.GIF http://www.astro.ulf.edu.aac/icons/solsyt.gif http://connect.ccsn.edu/edu/shs/grant/solar_system.gif http://www.bonus.com/bonus/card/solarsystembrowser/solarsystembrowaer.jpg

Acknowledgments

We would like to thank Shlomo Moran and Yoelle Maarek for useful discussions on the ideas presented in this paper.

REFERENCES

- ARIDOR, Y., CARMEL, D., LEMPEL, R., MAAREK, Y. S., AND SOFFER, A. 2000. Knowledge agents on the Web. *Proceedings of the 4th International Workshop on Cooperative Information Agents*.
- BHARAT, K. AND HENZINGER, M. R. 1998. Improved algorithms for topic distillation in a hyperlinked environment. *Proceedings of the 21st Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*.
- BRIN, S. AND PAGE, L. 1998. The anatomy of a large-scale hypertextual Web search engine. *Proceedings of the 7th International WWW Conference*.
- CASCIA, M. L., SETHI, S., AND SCLAROFF, S. 1998. Combining textual and visual cues for content-based image retrieval on the World Wide Web. *Proceedings of the IEEE Workshop on Content-based Access of Image and Video Libraries*.
- CHAKRABARTI, S., DOM, B., GIBSON, D., KLEINBERG, J., KUMAR, S., RAGHAVAN, P., RAJAGOPALAN, S., AND TOMKINS, A. 1999. Hypersearching the Web. *Scientific American*.
- CHAKRABARTI, S., DOM, B., GIBSON, D., KLEINBERG, J. M., RAGHAVAN, P., AND RAJAGOPALAN, S. 1998. Automatic resource list compilation by analyzing hyperlink structure and associated text. *Proceedings of the 7th International WWW Conference*.

¹⁸Although two of the URLs have a suffix of *.gif*, all three contain the same *.jpeg* image. The suffix, in these cases, does not describe the file type correctly.

- DEAN, J. AND HENZINGER, M. R. 1999. Finding related pages in the World Wide Web. *Proceedings of the 8th International World Wide Web Conference*.
- FLICKNER, M., SAWHNEY, H., NIBLACK, W., ASHLEY, J., HUANG, Q., DOM, B., GORKANI, M., HAFNER, J., LEE, D., PETKOVIC, D., STEELE, D., AND YANKE, P. 1995. The QBIC system. *IEEE Computer* 28, 9, IEEE Computer Society Press, Los Alamos, CA, 23–32.
- FRANKEL, C., SWAIN, M., AND ATHITSOS, V. 1996. Webseer: An image search engine for the World Wide Web. Tech. Rep. TR-96-14, Computer Science Dept., Univ. of Chicago.
- GEVERS, T. AND SMEULDERS, A. W. M. 1999. The PicToSeek www image search system. *Proceedings of the IEEE International Conference on Multimedia Computing and Systems* 1, 264–269.
- HARMANDAS, V., SANDERSON, M., AND DUNLOP, M. 1997. Image retrieval by hypertext links. *Proceedings of the 20th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, 296–303.
- KANTH, K., AGRAWAL, D., AND SINGH, A. 1998. Dimensionality reduction for similarity searching in dynamic databases. In *Proceedings of the ACM SIGMOD Int. Conf. on Management of Data* (1998), 166–176.
- KESSLER, M. 1963. Bibliographic coupling between scientific papers. *American Documentation* 14, 10–25.
- KLEINBERG, J. M. 1999. Authoritative sources in a hyperlinked environment. *J. ACM* 46: 5, 604–632.
- LEMPEL, R. AND MORAN, S. 2000. The stochastic approach for link-structure analysis (SALSA) and the TKC effect. *Proceedings of the 9th International WWW Conference*.
- LEMPEL, R. AND SOFFER, A. 2001. PicASHOW: Pictorial Authority Search by Hyperlinks on the Web. *Proceedings of the 10th International WWW Conference*.
- PENTLAND, A. P., PICARD, R. W., AND SCLAROFF, S. 1996. Photobook: Content-based manipulation of image databases. *Int. J. Compute. Vis.* 18, 3, 233–254.
- RAFIEL, D. AND MENDELZON, A. 2000. What is this page known for? computing Web page reputations. *Proceedings of the 9th International WWW Conference*.
- RUI, Y., HUANG, T., AND CHANG, S. 1999. Image retrieval: Past, present, and future. *J. Vis. Com. and Image Rep.* 10, 1–23.
- SMALL, H. 1973. Co-citation in the scientific literature: A new measure of the relationship between two documents. *J. American Soc. Info. Sci.* 24, 265–269.
- SMITH, J. R. AND CHANG, S.-F. 1996. Searching for images and videos on the World Wide Web. Tech. Rep. 459-96-25, Columbia University/CTR.
- VAN RIJSBERGEN, C. 1979. *Information Retrieval*. Butterworths.

Received April 2001; revised May 2001 and June 2001; accepted July 2001