

# Computing Geographical Scopes of Web Resources

Junyan Ding  
Computer Science Department  
Columbia University  
dingjy@cs.columbia.edu

Luis Gravano  
Computer Science Department  
Columbia University  
gravano@cs.columbia.edu

Narayanan Shivakumar  
Gigabeat, Inc.  
shiva@gigabeat.com

## Abstract

Many information resources on the web are relevant primarily to limited geographical communities. For instance, web sites containing information on restaurants, theaters, and apartment rentals are relevant primarily to web users in geographical proximity to these locations. In contrast, other information resources are relevant to a broader geographical community. For instance, an on-line newspaper may be relevant to users across the United States. Unfortunately, current web search engines largely ignore the *geographical scope* of web resources. In this paper, we introduce techniques for automatically computing the geographical scope of web resources, based on the textual content of the resources, as well as on the geographical distribution of hyperlinks to them. We report an extensive experimental evaluation of our strategies using real web data. Finally, we describe a geographically-aware search engine that we have built to showcase our techniques.

## 1 Introduction

The World-Wide Web provides uniform access to information available around the globe. Some web sites such as on-line stores and banking institutions are of “global” interest to web users world-wide, while many web sites contain information primarily of interest to web users in a geographical community, such as the

Bay Area or Palo Alto. Over the past few years, web users have been discovering web sites using web search engines such as AltaVista<sup>1</sup> and Google<sup>2</sup>. In practice, these engines are ineffective for identifying *geographically scoped* web pages. For instance, finding restaurants, theaters, and apartment rentals in or near specific regions is a difficult task with these web search engines.

Now consider the scenario in which we have a database with the geographical scope (e.g., a city, a state) of all “resources” (e.g., restaurants, newspapers) with a web presence. We can then exploit such information for a variety of applications, including the following:

- **Personalized searching:** Consider the case a resident in Palo Alto searches for “newspapers.” A geographically-aware search engine would first identify where the user is from (e.g., using a profile at [my.yahoo.com](http://my.yahoo.com) or [my.excite.com](http://my.excite.com)). The search engine then uses this information to return newspapers that are relevant to the user’s location, rather than returning references to newspapers all over the world. For instance, the engine might recommend The New York Times as a “globally relevant” newspaper, and the Stanford Daily as a local newspaper. Note that this strategy is not equivalent to the user querying the search engine for “newspaper AND Palo Alto,” since such a query would miss references to The New York Times, a newspaper that is published in a city not in the vicinity of Palo Alto. This newspaper even has the name of a specific city (“New York”) in its name, but is nevertheless geographically relevant to the entire United States.
- **Improved browsing:** Web portals like Yahoo! already classify web resources *manually* according to their geographical scope<sup>3</sup>. The techniques that we present in this paper will make it possible

---

*Permission to copy without fee all or part of this material is granted provided that the copies are not made or distributed for direct commercial advantage, the VLDB copyright notice and the title of the publication and its date appear, and notice is given that copying is by permission of the Very Large Data Base Endowment. To copy otherwise, or to republish, requires a fee and/or special permission from the Endowment.*

<sup>1</sup><http://www.altavista.com>

<sup>2</sup><http://www.google.com>

<sup>3</sup><http://info.ford.osar.com/Regional/>

to conduct such hierarchical categorization efforts automatically, improving their scalability.

It is easy to build geographically aware applications such as the above if we are supplied with a table that lists the geographical scope of each resource. Unfortunately, no such table exists for web resources. In this paper, we consider how to mine the web and automatically construct such a table using web hyperlinks and the actual content of web pages. For example, we can map every web page to a location based on where its hosting site resides. Then, we can consider the location of all the pages that point to, say, the Stanford Daily home page<sup>4</sup>. By examining the distribution of these pointers we can conclude that the Stanford Daily is of interest mainly to residents of the Stanford area, while The Wall Street Journal is of nation-wide interest. We can draw the same conclusion by analyzing the geographical locations that are mentioned in the pages of the Stanford Daily and in those of The Wall Street Journal.

The primary contributions of this paper include:

1. **Algorithms to estimate geographical scope:** We propose a variety of algorithms that automatically estimate the geographical scope of resources, based on exploiting either the distribution of HTML links to the resources (Section 3) or the textual content of the resources (Section 4).
2. **Measures to evaluate quality of algorithms:** We introduce evaluation criteria for our estimation algorithms, based on traditional information-retrieval metrics (Section 5).
3. **Experimental study of techniques:** We empirically evaluate our algorithms using real web data (Section 6).
4. **Implementation of a geographically aware search engine:** We also discuss how we used our algorithms in the implementation of a geographically aware search engine for on-line newspapers, which is accessible at <http://www.cs.columbia.edu/~gravano/GeoSearch> (Section 7).

## Related Work

Traditional information-retrieval research has studied how to best answer keyword-based queries over collections of text documents [14]. These collections are typically assumed to be relatively uniform in terms of, say, their quality and scope. With the advent of the web, researchers are studying other “dimensions” to the data that help separate useful resources from

less-useful ones in an extremely heterogeneous environment like the web. Techniques for text-database selection [3, 8, 13, 10] decide what web databases to use to answer a user query, basing this decision on the textual contents of the web databases.

Recent research has started to exploit web links for improving web-page categorization [4] and for web mining [6, 9, 5]. Notably, search engines such as Google [1] and HITS [6, 12] estimate the “importance” of web pages by considering the number of hyperlinks that point to them. The rationale for their heuristics is that the larger the number of web users who made a hyperlink to a web page, the higher must be the importance of the page. In essence, this work manages to capture an additional dimension to the web data, namely how important or authoritative the pages are. Unlike the new techniques that we introduce in this paper, HITS and Google ignore the *spatial distribution* of incoming links to a web resource.

In this paper, we propose to extract yet another crucial dimension of the web data, namely the geographical scope of web resources. This new dimension can then be used to complement traditional information retrieval techniques and those used by Google and HITS to answer web queries in more effective ways. Some commercial web sites already *manually* classify web resources by their location, or keep directory information that lists where each company or web site is located (e.g., see <http://www.iatlas.com>). Quite recently, the NorthernLight search engine<sup>5</sup> has started to extract addresses from web pages, letting users narrow their searches to specific geographical regions (e.g., to pages “originated” within a five-mile radius of a given zip code). Users benefit from this information because they can further filter their query results. In reference [2], we discussed how to map a web site (e.g., <http://www-db.stanford.edu>) to a geographical location (e.g., Palo Alto), and we also presented a tool to visualize such geographical web data. In this paper, we extend this preliminary work to a harder problem: how to *automatically* estimate the geographical *scope* of a web resource? That is, which data is targeted towards residents of a city as opposed to the country, or the world?

## 2 Geographical Scopes of Web Resources

Web resources are built with a target audience in mind. Sometimes this audience is geographically enclosed in some neighborhood (e.g., the target audience of the web page of a local pizzeria that delivers orders to houses up to 2 miles away from the store). Some other times, the target audience of a resource is distributed across the country (e.g., the target audience of the web page of the USA Today newspaper). In this section,

<sup>4</sup>Citations from pages hosted on national access providers like America On Line would be ignored in this process, unless we can map these citations to the physical location of their creator. We discuss this issue further in Section 6.1.

<sup>5</sup><http://www.northernlight.com/geosearch.html>

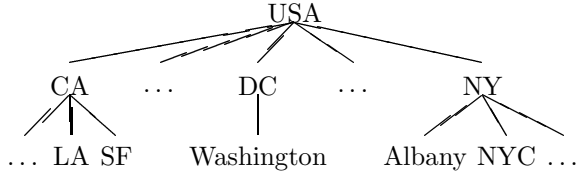


Figure 1: Portion of the hierarchy of geographical locations for the United States.

we introduce the notion of the *geographical scope* of a web resource, which captures the geographical distribution of the target audience of a web resource. This notion is a subjective one, the same way that the information-retrieval notion of *document relevance* is subjective [14].

**Definition 1:** The *geographical scope* of a web resource  $w$  is the geographical area that the creator of  $w$  intends to reach.

Using this informal definition, the geographical scope of our pizzeria above is the neighborhood where the pizzeria resides, whereas the geographical scope of the USA Today newspaper is the entire United States.

For concreteness, in the rest of the paper we focus on how to approximate the geographical scope of web resources within the United States. For this, we will view the United States as a three-level location hierarchy (Figure 1). The root of the hierarchy corresponds to the entire country. The next level down the hierarchy has one node for each of the 50 states, plus one node for the District of Columbia. Finally, the leaf level in the hierarchy has one node for each city in the country. Using this hierarchy, a human expert might specify that the geographical scope of the USA Today newspaper is the whole USA. In contrast, the geographical scope of The Arizona Daily Star Online is the state of Arizona, since this state is the target audience of this newspaper. Finally, the geographical scope of yet another newspaper, The Knoxville News-Sentinel, has the city of Knoxville as its geographical scope. Of course, this three-level hierarchy can be extended to span all the countries in the world, as well as to further localize resources in cities, to counties and boroughs. However, for simplicity this paper focuses only on the three levels listed above.

Given our three-level hierarchy of geographical locations in the United States, we can choose to define the geographical scope of web resources in different ways. For example, instead of indicating that the geographical scope of the USA Today newspaper is the top node in the hierarchy (i.e., the whole United States), we could list all 50 states plus Washington D.C. as comprising this geographical scope. Although it could be argued that this state-level formulation expresses the same information as the country-level one, we will

always express geographical scopes using nodes that are “as high” as possible in our three-level hierarchy. Thus, instead of aggregating the information that the USA Today is a national newspaper out of the list of states of its geographical scope, we will state this fact directly, and simply specify its scope to be the United States as a whole.

As mentioned above, the notion of geographical scope is subjective. To capture this notion accurately, we could hand-classify each web resource according to its intended geographical scope. (Incidentally, this is the way that web portals like Yahoo! operate.) In this paper, we study scalable ways to *automatically* approximate the resources’ geographical scopes. Sections 3 and 4 describe two ways in which we can estimate the geographical scope of a web resource. Later, Section 6 will report the experiments that show that our automatically-computed approximations closely match the “ideal,” subjective definition.

### 3 Exploiting Resource Usage

In this section, we show how we can estimate the geographical scope of web resources by exploiting the link structure of the web. (We will present an alternative estimation method that exploits the *contents* of the web resources in Section 4.)

Consider a web resource whose geographical scope is the entire United States (e.g., the USA Today newspaper). Such a resource is likely to then attract interest across the country. Our assumption in this section is that this interest will translate in web pages across the country containing HTML links to this web resource.<sup>6</sup> Conversely, a resource with a much more limited geographical scope will exhibit a significantly different link distribution pattern across the country. Hence a promising way to estimate the geographical scope of a resource is to study the geographical distribution of links to the resource. More specifically, two conditions that a location  $\ell$  will have to satisfy to be in the geographical scope of a resource  $w$  are:

- A significant *fraction* of  $\ell$ ’s web pages contain links to  $w$  (Section 3.1).
- The web pages in  $\ell$  that contain links to  $w$  are distributed *smoothly* across  $\ell$  (Section 3.2).

Below we show how to estimate the geographical scope of a web resource  $w$  by identifying a set of candidate locations  $\ell$  that satisfy the two conditions above. This process results in the estimated geographical scope of  $w$ . Our experiments of Section 6 will show that these

<sup>6</sup>Of course, if we knew who *accesses* each web resource we could use this information for our problem. Unfortunately, web access logs for all resources whose geographical scopes we would like to characterize are not easily available.

estimates are often a good approximation of the subjective geographical scopes that we discussed in Section 2.

### 3.1 Measuring Interest: *Power*

Intuitively, a location  $\ell$  that is in the geographical scope of a web resource  $w$  should exhibit relatively high “interest” in  $w$  among its web pages. In other words, a relatively high fraction of the web pages originated in  $\ell$  should contain links to resource  $w$ .  $Power(w, \ell)$  measures the relative interest in  $w$  among the pages in location  $\ell$ :

$$Power(w, \ell) = \frac{Links(w, \ell)}{Pages(\ell)} \quad (1)$$

where  $Links(w, \ell)$  is the number of pages in location  $\ell$  that contain a link to web resource  $w$ , and  $Pages(\ell)$  is the total number of web pages in  $\ell$ . (We explain how we compute these numbers in Section 6.1.)

### 3.2 Measuring Uniformity: *Spread*

Geographical locations can be decomposed into several sub-locations. As an example, the United States consists of 50 states plus the District of Columbia, while the state of New York, in turn, comprises a number of cities (e.g., Albany, New York City). As we discussed in the previous section, to include a location  $\ell$  (e.g., the state of New York) in the geographical scope of a web resource  $w$ , there should be a sufficiently high “interest” in resource  $w$  in location  $\ell$  (i.e.,  $Power(w, \ell)$  is high). In addition, we need to ask that this interest be spread smoothly across the location. Thus, a resource with an unusually high number of links originating, say, in New York City, but with no links coming from other New York state cities should not have the state of New York in its geographical scope, but perhaps just New York City instead.

To determine how uniform the distribution of links to a web resource  $w$  is across a location  $\ell$ , we introduce a second metric, *Spread*. Intuitively,  $Spread(w, \ell)$  will be high whenever  $Power(w, \ell_i) \sim Power(w, \ell_j)$  for all “sub-locations”  $\ell_i, \ell_j$  that are children of  $\ell$  in the location hierarchy of Section 2. In what follows, we provide three alternative definitions of *Spread*. These definitions are all built on this intuition, but will compute the value of  $Spread(w, \ell)$  using techniques borrowed from different fields. In Section 6 we experimentally compare how these three definitions perform relative to each other.

For our three definitions of *Spread*,  $Spread(w, \ell)$  will have the maximum possible value (i.e., a value of 1) in the following two special cases:

- $\ell$  is a leaf node of our location hierarchy: In this case, by definition, the distribution of *Power* across  $\ell$  is completely uniform, because we regard

$\ell$  as an “atomic” location. In this paper, these atomic locations are the United States cities.

- $Power(w, \ell)=0$ : In this case, there is no “interest” at all in resource  $w$  across location  $\ell$ . Since *Spread* measures the uniformity of this interest across  $\ell$ ,  $Spread(w, \ell)$  is trivially maximum in this case.

Next, we give three alternative definitions for  $Spread(w, \ell)$  for the case when  $\ell$  is not a leaf node in our location hierarchy and  $Power(w, \ell) > 0$ . In the definitions below,  $\ell_1, \dots, \ell_n$  are the children of  $\ell$  in the hierarchy. Also, we associate with location  $\ell$  vector  $\vec{Pages} = (p_1, \dots, p_n)$ , which lists the number of pages  $p_i = Pages(\ell_i)$  of each child  $\ell_i$  of  $\ell$ . A second vector associated with  $\ell$ ,  $\vec{Links} = (l_1, \dots, l_n)$ , lists the number of pages  $l_i = Links(w, \ell_i)$  that have a link to resource  $w$  at location  $\ell_i$ , for  $i = 1, \dots, n$ . Finally, vector  $\vec{Power} = (r_1, \dots, r_n)$  lists the value of  $Power$   $r_i = Power(w, \ell_i)$  for each sub-location of  $\ell$ .

#### Vector-Space Definition of *Spread*

The first definition of *Spread* is inspired in the vector-space model from information retrieval [14]. Intuitively, we will compute how “similar” vectors  $\vec{Pages}$  and  $\vec{Links}$  are by computing the cosine of the angle between them. If the fraction of pages with links to  $w$  is mostly constant across all of  $\ell$ ’s children  $\ell_1, \dots, \ell_n$ , then  $\vec{Pages}$  and  $\vec{Links}$  will be roughly scaled versions of one another, and the cosine of the angle between these vectors will be close to 1:

$$\begin{aligned} Spread(w, \ell) &= \vec{Pages} \odot \vec{Links} \\ &= \frac{\sum_{i=1}^n p_i \times l_i}{\sqrt{\sum_{i=1}^n p_i^2} \cdot \sqrt{\sum_{i=1}^n l_i^2}} \end{aligned} \quad (2)$$

#### Relative-Error Definition of *Spread*

Let  $R = (\sum_{i=1}^n l_i) / (\sum_{i=1}^n p_i)$ . If the distribution of interest in  $w$  were perfectly smooth, then  $r_i = R$  for all  $i$ . To measure how far we are from this perfectly smooth distribution, we compute how much each  $r_i$  deviates from the “target” value  $R$ . We can then give a definition of *Spread* based on computing the “relative error” for each  $\ell_i$  with respect to  $R$ :

$$Spread(w, \ell) = \frac{1}{1 + \sum_{i=1}^n \frac{1}{p_i} \sum_{i=1}^n p_i \cdot \frac{|R - r_i|}{R}} \quad (3)$$

#### Entropy Definition of *Spread*

Our third and final definition for *Spread* is based on the notion of entropy from information theory [11]. To give this definition, we assume that there is an “information source” associated with web resource  $w$  and geographical location  $\ell$ . The information source generates symbols representing the different children of

$\ell$ , namely  $\ell_1, \dots, \ell_n$ . Moreover, we assume that this information source generates its symbols by infinitely executing three steps:

1. Randomly select an  $\ell_i$ .
2. Randomly select a web page located in  $\ell_i$ .
3. If the web page has a link to web site  $w$ , then generate a symbol representing  $\ell_i$ .

Intuitively, when  $r_i = \text{Power}(w, \ell_i)$  is uniform across the  $\ell_i$  sub-locations, the information source will achieve the maximum entropy available at geographical location  $\ell$ , which is  $\log n$ . To make this definition comparable across geographical locations with different numbers of sub-locations, we define *Spread* as follows:

$$\text{Spread}(w, \ell) = \frac{-\sum_{i=1}^n \frac{r_i}{\sum_{j=1}^n r_j} \cdot \log(\frac{r_i}{\sum_{j=1}^n r_j})}{\log n} \quad (4)$$

### 3.3 Estimating Geographical Scopes

The previous sections showed metrics to measure the strength ( $\text{Power}(w, \ell)$ ) and uniformity ( $\text{Spread}(w, \ell)$ ) of the interest in a web resource  $w$  at a location  $\ell$ . In this section we define how we can use *Power* and *Spread* to estimate what locations we should include in the geographical scope of a given web resource.

As a first step to estimate the geographical scope of a web resource  $w$ , we identify the locations  $\ell$  in our hierarchy of Section 2 with  $\text{Spread}(w, \ell) \geq \tau_c$ , for some given threshold  $0 \leq \tau_c \leq 1$ . These are the locations with a relatively smooth distribution of links to  $w$  across their sub-locations. Furthermore, we only include in  $\text{CGS}(w)$ , the *candidate geographical scope* for  $w$ , those locations that have no ancestor  $\ell'$  with  $\text{Spread}(w, \ell') \geq \tau_c$ . In other words,  $\text{CGS}(w)$  contains locations with smooth distribution of links for  $w$  such that are not “subsumed” by any other ancestor location also in  $\text{CGS}(w)$ :

**Definition 2:** *The candidate geographical scope  $\text{CGS}(w)$  of a web resource  $w$  is a set of nodes in the geographical hierarchy. A location  $\ell$  is in  $\text{CGS}(w)$  if it satisfies the following two conditions, given a fixed threshold  $\tau_c$ :*

- $\text{Spread}(w, \ell) \geq \tau_c$ .
- For all  $\ell'$  that is an ancestor of  $\ell$ ,  $\text{Spread}(w, \ell') < \tau_c$ .

Given a web resource  $w$ , we can compute  $\text{CGS}(w)$  with a simple algorithm that recursively visits the nodes in the location hierarchy top-down.<sup>7</sup>

<sup>7</sup>We have investigated an alternative, “stricter” definition of  $\text{CGS}(w)$ . According to this definition,  $\ell \in \text{CGS}(w)$  if every location  $\ell'$  in the location subtree rooted at  $\ell$  has  $\text{Spread}(w, \ell') \geq \tau_c$ . Our experimental results showed that the weaker definition that we give above outperformed this stricter definition. For space constraints, we then do not discuss this stricter version further.

The candidate geographical scope of a resource  $w$ ,  $\text{CGS}(w)$ , contains locations exhibiting relatively smooth interest in  $w$ . However, as we discussed earlier, this interest could be quite small in some cases. In particular, a location  $\ell$  with  $\text{Power}(w, \ell)=0$  (e.g., a leaf node) might be included in  $\text{CGS}(w)$ , which is clearly undesirable. Consequently, we need to prune our candidate geographical scopes to only include locations with high enough *Power* in the final estimated geographical scope of a resource:

**Definition 3:** *The estimated geographical scope  $\text{EGS}(w)$  of a web resource  $w$  is a set of locations obtained from  $\text{CGS}(w)$  using one of the following scope pruning strategies:*

- **Top- $k$  pruning:** *Given an integer  $k$ ,  $\text{EGS}(w)$  consists of the top- $k$  locations in  $\text{CGS}(w)$ , in decreasing order of their *Power*.*
- **Absolute-threshold pruning:** *Given a threshold  $\tau_e$ ,  $\text{EGS}(w) = \{\ell \in \text{CGS}(w) | \text{Power}(w, \ell) \geq \tau_e\}$ .*
- **Relative-threshold pruning:** *Given a percentage  $p$ ,  $\text{EGS}(w) = \{\ell \in \text{CGS}(w) | \text{Power}(w, \ell) \geq \text{max\_Power}(w) \times p\}$ , where  $\text{max\_Power}(w) = \max\{\text{Power}(w, \ell) | \ell \in \text{CGS}(w)\}$ .*

## 4 Exploiting Resource Contents

So far, we have used the distribution of links to a resource to estimate the resource’s geographical scope. A natural question, however, is whether we can instead just examine the resource’s contents to accomplish this task. In this section we explore this idea, and discuss how to use the resources’ text to estimate their geographical scope.

Consider a resource whose geographical scope is, say, the state of New York. We may argue that the text in such a resource is likely to mention New York cities more frequently than locations corresponding to other states or countries. This is our main assumption in this section. (Section 6 experimentally compares the resulting technique with our link-based strategy of Section 3.) Hence an interesting direction to explore to estimate the geographical scope of a resource is to study the distribution of locations that are *mentioned* in the resource. More specifically, two conditions that a location  $\ell$  will have to satisfy to be in the geographical scope of a resource  $w$  are:

- A significant *fraction* of all locations mentioned in  $w$  are either  $\ell$  itself or a sub-location of  $\ell$ .
- The location references in  $w$  are distributed *smoothly* across  $\ell$ .

Next, Section 4.1 shows that we can use the location references in the contents of a web resource to define

a variation of the *Power* and *Spread* metrics of Section 3. We then estimate the geographical scopes completely analogously as we did for the link-based strategy. Later, Section 4.2 addresses a fundamental step in our content-based approach, namely how we can effectively extract the location names from the text of a resource.

#### 4.1 Estimating Geographical Scopes

To estimate whether a location  $\ell$  is part of a resource  $w$ 's geographical scope we will proceed exactly as in Section 3 and compute (modified versions of)  $Power(w, \ell)$  and  $Spread(w, \ell)$ . For this, we need to extract from  $w$  two numbers. The first one,  $Locations(w)$ , is the number of references to geographical locations in  $w$ 's text. The second one,  $References(w, \ell)$ , is the number of references to  $\ell$  mentioned in  $w$ 's text.<sup>8</sup> Given these counts, we can adapt our definition of *Power* from Section 3.1 in the following way:

$$Power(w, \ell) = \frac{References(w, \ell)}{Locations(w)} \quad (5)$$

To adapt the definition of *Spread* of Section 3.2, we now define the following three vectors for a web resource  $w$  and a location  $\ell$  with children  $\ell_1, \dots, \ell_n$ . First, vector  $\vec{Locations} = (p_1, \dots, p_n)$  is a vector with every element having the same value  $p_i = Locations(w)$ , which is the number of references to geographical locations in  $w$ 's text. Second, vector  $\vec{References} = (l_1, \dots, l_n)$  lists the number of references to each sub-location  $\ell_i$  in  $w$ 's text, i.e.,  $l_i = References(w, \ell_i)$ . Finally, vector  $\vec{Power} = (r_1, \dots, r_n)$  lists each sub-location's *Power* value  $r_i = Power(w, \ell_i)$ . These vectors will play a role that is completely analogous to those of the  $\vec{Pages}$ ,  $\vec{Links}$ , and  $\vec{Power}$  vectors of Section 3, respectively, for defining  $Spread(w, \ell)$ . We can now use exactly the same definitions for *Spread* that we used in Section 3 and calculate the estimated geographical scope  $EGS(w)$  for a web resource  $w$ .

#### 4.2 Extracting and Processing Location References

To estimate the geographical scope of a web resource  $w$  as in the previous section, we need to extract all of the locations that are mentioned in the textual contents of  $w$ . Furthermore, the technique above expects the list of *cities* that are mentioned in the text of the web resources. In this section, we discuss the main problems involved in such an extraction process.

<sup>8</sup>We will discuss in Section 4.2 how we map references to, say, an entire state to references to individual cities within the state, which is what we count in *References* and *Locations*.

#### Extracting Location Names from Plain Text

State-of-the-art named-entity taggers manage to identify entities like people, organizations, and locations in natural-language text with high accuracy. For the experiments that we report in Section 6 we used the Alembic Workbench system developed at MITRE [7].

#### Normalizing and Disambiguating Location Names

After the tagging phase in which we identify the locations (e.g., "New York City," "California") mentioned in  $w$ , we should map each location to an unambiguous city-state pair. Problems that arise when completing this task include:

- **Aliasing:** Different names might be commonly used for the same location. For example, San Francisco is often referred to as SF. It is relatively easy to address this problem at the country or state level. (These aliases are indeed quite limited, and we compiled a list of them by hand.) For cities, though, we resorted to a web-accessible database of the United States Postal Service (USPS)<sup>9</sup>. For each zip code, this service returns a list of variations of the corresponding city's name. For example, if we use Columbia University's zip code, 10027, we obtain a list of names for New York City, including New York, Manhattan, New York City, NY City, NYC, and, interestingly enough, Manhattanville. (Incidentally, the USPS standard form for this city is New York.) By repeatedly querying the USPS database with different zip codes, we can build a list of city-name aliases, together with the corresponding "normal form" for each group.
- **Ambiguity:** Another problem when processing a given city name is that it can refer to cities in different states. For example, four states, Georgia, Illinois, Mississippi, and Ohio, have a city called Columbus. A reference to such a city without a state qualification is inherently ambiguous, unless of course we could understand the context in which the reference was made. We have developed heuristics for managing this kind of ambiguous location references. Our technique starts by identifying the unambiguous location references in the web resource at hand  $w$ , and uses them to disambiguate the remaining references. Intuitively, if  $w$  mentions mostly locations in the state of New York, for example, we will assume that a reference to "Manhattan" is a reference to New York City, not to Manhattan, Kansas. More specifically, if  $w$  mentions an ambiguous city name  $C$   $m$  times, and  $C$  can refer to a city in a number of

<sup>9</sup><http://www.usps.gov>

states  $S_1, \dots, S_k$ , then we “distribute” the  $m$  occurrences of  $C$  among the  $k$  states proportionally to the distribution of unambiguous cities in these states. Suppose that in our example 90% of the unambiguous cities that are mentioned in resource  $w$  are in the state of New York, and the remaining 10% are in the state of Kansas. Then, if  $w$  refers to Manhattan five times, we will assume that 4.5 of these references correspond to New York, NY, and only 0.5 of them to Manhattan, KS.

## Mapping Locations to City Names

A location name can refer to a city, a state, or a country, for example. Our technique to estimate geographical scopes analyzes the distribution of *cities* that are mentioned at a web resource  $w$ . Consequently, we need a way to map references to, say, states to city references that our technique can use. For this, we simply “push down” references to high-level locations in our location hierarchy (Section 2). This way, a reference to the state of New York will be pushed down as a reference to every city in the state. When we propagate these references down, we also scale their *weight* by some constant  $\alpha$ . (A value of  $\alpha = 0.1$  worked best in our Section 6 experiments.)

## 5 Evaluating the Quality of the Estimated Geographical Scopes

In the previous sections we discussed two approaches to estimating the geographical scope of resources. Of course, other approaches are possible (e.g., a “hybrid” strategy combining our two techniques). We now propose measures to evaluate the quality of any such algorithm for estimating a web resource’s geographical scope.

To evaluate the quality of our estimated geographical scopes, we need to compare them against the ideal, subjective scopes. We could base our comparison on metrics commonly used for classification tasks: for example, we could just compute the number of web resources in our testbed for which we managed to identify their geographical scope *perfectly*. Such a metric would not fully capture the nuances of our problem. For example, if the geographical scope of a resource  $w$  is {California} and we compute  $EGS(w)$  as, say, {California, New York City}, this metric would mark our answer as completely wrong. Similarly, consider the case where our  $EGS(w)$  computation consists of, say, 90% of the California cities, but does not include California as a whole state, which would have been the perfect answer. Traditional classification accuracy metrics would also consider our estimate as completely wrong, even when our technique managed to identify only cities in the right state as part of the geographical scope of  $w$ .

With these observations in mind, we adapt the precision and recall metrics from the information retrieval field to yield metrics that we believe are appropriate for our problem. More specifically, we will define *precision* and *recall* for our problem as follows, after we introduce an auxiliary definition. Given a set of locations  $L$ , we will “expand” it by including all locations under a location  $\ell \in L$ . Thus,  $Expanded(L) = \{\ell' \text{ location} \mid \ell' \in L \text{ or } \ell' \text{ is in the location subtree of some } \ell \in L\}$ . Now, let  $w$  be a web resource, *Ideal* be its “expanded” geographical scope, and *Estimated* be our expanded estimate,  $Expanded(EGS(w))$ . Then:

$$\begin{aligned} Precision(w) &= \frac{|Ideal \cap Estimated|}{|Estimated|} \\ Recall(w) &= \frac{|Ideal \cap Estimated|}{|Ideal|} \end{aligned}$$

Intuitively, precision measures the fraction of locations in an estimated geographical scope that are correct, i.e., that are also part of the ideal geographical scope. (Perfect precision might be trivially achieved by always returning empty geographical scopes.) Recall measures the fraction of the locations in the ideal geographical scope that are captured in our estimated geographical scope. (Perfect recall might be trivially achieved by always including all locations in the geographical scopes.) Finally, to simplify the interpretation of our experiments, we combine precision and recall into a single metric using the *F-measure* [15]:

$$F(w) = \frac{2 \times Precision(w) \times Recall(w)}{Precision(w) + Recall(w)}$$

## 6 Experimental Evaluation

Section 2 defined the “ideal,” subjective geographical scope of a web resource  $w$ . Later, we showed how we can automatically calculate the estimated geographical scope  $EGS(w)$  by analyzing the geographical distribution of HTML links to  $w$  (Section 3), or, alternatively, by analyzing the distribution of location names from the textual contents of  $w$  (Section 4). In this section, we experimentally evaluate how well our different techniques can approximate the ideal geographical scopes using the evaluation criteria we discussed in Section 5. We describe our experimental setting in Section 6.1. We then report the results of our experiments, which involved real web resources, in Section 6.2.

### 6.1 Experimental Setting

In this section we explain the main aspects of our experimental setting. In particular, we describe the real web resources that we used, and highlight some of the challenging implementation issues that we had to address to carry out our study.

## Web Resources

Ideally, to evaluate our techniques of Sections 3 and 4 we should use a set of real web resources, each with its corresponding geographical scope, as determined by a human expert. (Analogously, the information retrieval field relies on human relevance judgments to evaluate the performance of search-engine algorithms [14].) For our experiments, we needed a list of web resources whose intended geographical scope was self-apparent and uncontested. Furthermore, we wanted our list to cover the different levels of our location hierarchy of Section 2. In other words, we wanted resources who would have the United States as their geographical scope, but we also wanted resources whose geographical scope was at the state and city levels. Finally, the resources that we picked needed to have a sufficiently large number of HTML links directed to them, so that we can apply our technique of Section 3. (We discuss how to handle resources with not enough references to them in Section 8.) With the above goals in mind, we collected a list of 150 web resources whose geographical scopes span the three levels of our location hierarchy:

- **National level:** 50 of our web resources have the United States as their geographical scope. These resources are the 50 most heavily cited Federal Government web sites listed in the FedWorld web site <sup>10</sup>. (We determined the 50 most cited pages by querying AltaVista to obtain the number of pages with links to each of these resources.) These web sites have the whole United States as their intended audience, and include the web sites of NASA <sup>11</sup> and the National Endowment for the Arts <sup>12</sup>, for example.
- **State level:** 50 of our web resources have a state as their geographical scope. These resources are the official web site of each state in the United States (e.g., <http://www.state.ny.us> (state of New York)), and have their corresponding state as their geographical scope.
- **City level:** 50 of our web resources have a city as their geographical scope. These resources are the 50 most cited among the US cities' official web sites (e.g., <http://www.ci.sf.ca.us> (San Francisco)), and have their corresponding city as their geographical scope. (We obtained a list of the US cities' official web sites from Piper Resources' "State and Local Government on the Net." <sup>13</sup>)

## Implementation Issues

We now describe some interesting tasks that we had to perform to run our experiments:

- **Mapping web pages to city names:** Our technique of Section 3 requires that we find all pages with HTML links to a given web resource  $w$ . After identifying these pages, we need to place their location so that we can study their geographical distribution and estimate  $w$ 's geographical scope. This is a challenging task, because we really need the location of the *author* of a page, which might be quite different from the location of the site that hosts the page. For example, web pages with links to  $w$  from, say, the [aol.com](http://aol.com) domain are hardly useful for our task: If we examine just the location of the web site where these pages reside, we would most likely be misguided in determining  $w$ 's geographic scope. (We elaborated on these issues further and outlined alternative approaches to "placing web pages on the map" in [2].) A key observation that we exploit for our experiments is that it suffices for our Section 3 technique to have a reasonable *sample* of the pages with links to resource  $w$  to estimate  $w$ 's geographical scope. Following this observation, we focussed on web pages whose author's location we could determine reliably, and that would span the entire United States. More specifically, we analyzed link information from pages originating only in educational domains (e.g., from web sites with a [.edu](http://.edu) suffix, like [www.columbia.edu](http://www.columbia.edu)). Given such a page, we query the `whois` service to map the page's web site into its corresponding zip code. After this, we query the USPS zip-code server and obtain the standard city name associated with the zip code.
- **Refining our location hierarchy:** Our experimental setting considers links originating only in educational institutions. Unfortunately, not all cities have one such institution. Hence, we refined our location hierarchy to include only cities with a university with a [.edu](http://.edu) web site. We further pruned our list by eliminating every city hosting fewer than 500 pages in [.edu](http://.edu) web sites, so that we analyze only cities with a significant web presence in [.edu](http://.edu) domains. At the end of this process, we are left with a location hierarchy consisting of 673 cities as leaf nodes, the 50 states and the District of Columbia as intermediate nodes, and the entire United States as the root node.
- **Computing  $Pages(\ell)$  and  $Links(w, \ell)$ :** For each city  $\ell$  in our location hierarchy, we need to obtain the number of pages from [.edu](http://.edu) domains that are located in it. To get this number, we query AltaVista and obtain the number of pages that each educational institution in location  $\ell$  hosts. By adding these numbers for each institution in  $\ell$  we compute  $Pages(\ell)$ , which we need in Section 3. Similarly, we can identify how many of these pages have links to a specific web resource

<sup>10</sup><http://www.fedworld.gov/locator.htm>

<sup>11</sup><http://www.nasa.gov>

<sup>12</sup><http://www.arts.endow.gov>

<sup>13</sup><http://www.piperinfo.com/state/states.html>



$w$  to compute  $Links(w, \ell)$ .

- **Obtaining the textual contents of a web resource  $w$ :** For our content-based technique of Section 4, we download the full-text contents of the 150 web sites of our testbed, using Gnu’s `wget` web crawler. We then use the `lynx` browser to eliminate the HTML tags in the web pages and extract the plain-English text in them. As explained in Section 4, after this we run the Alembic named-entity tagger [7] to extract the location names that are mentioned in the plain text. Finally, we resolve aliasing and ambiguity issues, map locations to city names, and estimate the geographical scopes as outlined in Section 4.<sup>14</sup>

## 6.2 Experimental Results

In Table 1 we summarize the algorithms we now evaluate. In addition to the three different definitions of *Spread* we discussed in Section 3.2, we also consider the following two simple “baseline” algorithms for computing candidate geographical scopes. The first one, *AllLeaves*, always defines the candidate geographical scope  $CGS(w)$  of a web resource  $w$  as consisting of all of the cities in our location hierarchy. In contrast, the second baseline technique, *OnlyRoot*, always defines  $CGS(w)$  as consisting of the United States only. The candidate geographical scope, obtained by any baseline technique or *Spread* definition, is then pruned by one of the three scope-pruning strategies to produce the estimated geographical scope as described in Section 3.3. Table 1 also summarizes the parameters involved with each of the different algorithms. For example, recall that  $k$  is a tunable parameter in the *TopK* scope-pruning strategy of Section 3.3.

We comprehensively evaluated the above algorithms to understand the impact of the different tunable parameters on precision, recall, and the  $F$ -measure. Due to lack of space, we present a few sample results to highlight some of our key observations. Specifically, we present our results for the relative-threshold pruning strategy *RelThr*. We evaluated our results for the *TopK* and for the *AbsThr* strategies as well, and observed similar trends. Hence we do not discuss these further.

In Figure 2, we show the impact of parameter  $p$  on the average  $F$ -measure for the link-based approach using the relative-threshold pruning strategy *RelThr*. (We use the values of  $\tau_c$  that are specified in Table 2.) Notice that all the *Spread* definitions perform very well, especially as  $p$  increases, and the *Spread* definitions have a much higher average  $F$ -measure compared to the strawman *AllLeaves* and *OnlyRoot* tech-

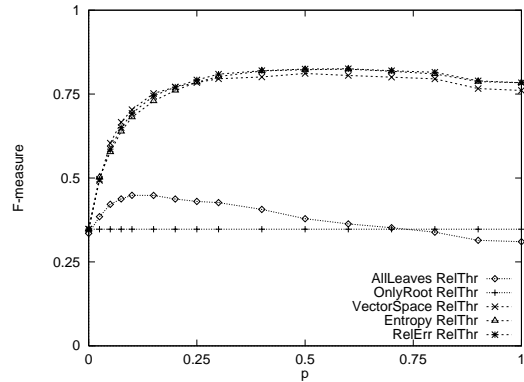


Figure 2: Average  $F$ -measure for the link-based strategy of Section 3 as a function of  $p$  (*RelThr* pruning strategy).

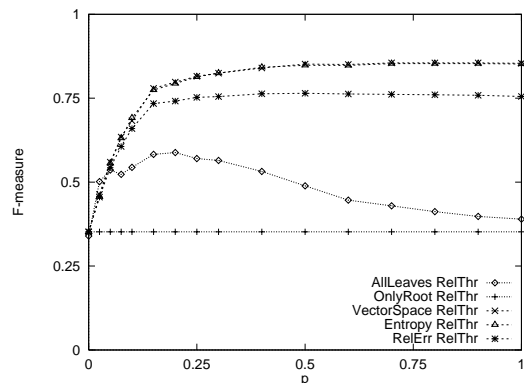


Figure 4: Average  $F$ -measure for the content-based strategy of Section 4 as a function of  $p$  (*RelThr* pruning strategy).

niques. Figure 3 shows the average precision and recall for *RelErr*, with the *RelThr* pruning strategy. Our techniques have more than 75% average precision and recall for several settings of the  $p$  parameter, which translates into the correspondingly high average  $F$ -measure values of Figure 2.

So far we have evaluated our link-based techniques for different parameter settings. We now discuss similar results for our content-based techniques of Section 4. In Figure 4 we report the impact of  $p$  on the average  $F$ -measure for the content-based approach on our entire data set, using the values of  $\tau_c$  specified in Table 3. We observe similar results to our link-based approach in that our techniques have high average  $F$  values, especially for  $p > 0.2$ , compared to the strawman techniques.

Table 2 summarizes our results from the previous graphs for the link-based approach, and reports the “best” (i.e., highest average  $F$  value) parameter values

<sup>14</sup>We only used 142 of the 150 web resources in our testbed to evaluate the content-based technique of Section 4: out of the remaining eight web resources, either `wget` could not crawl their pages, or the named-entity tagger that we used, Alembic, could not find any location name in their pages.

	<i>Label</i>	<i>Description</i>	<i>Associated Parameter</i>
<i>Baseline Techniques</i>	<i>AllLeaves</i>	Scope consists of all USA cities	—
	<i>OnlyRoot</i>	Scope consists of just USA	—
<i>Spread Definition</i> (Section 3.2)	<i>VectorSpace</i>	Vector-space definition of <i>Spread</i>	$\tau_c$
	<i>Entropy</i>	Entropy definition of <i>Spread</i>	$\tau_c$
	<i>RelErr</i>	Relative-error definition of <i>Spread</i>	$\tau_c$
<i>Scope-Pruning Strategies</i> (Section 3.3)	<i>TopK</i>	Top- $k$ pruning	$k$
	<i>AbsThr</i>	Absolute-threshold pruning	$\tau_e$
	<i>RelThr</i>	Relative-threshold pruning	$p$

Table 1: The variations of our techniques that we use in our experiments, together with their associated parameters.

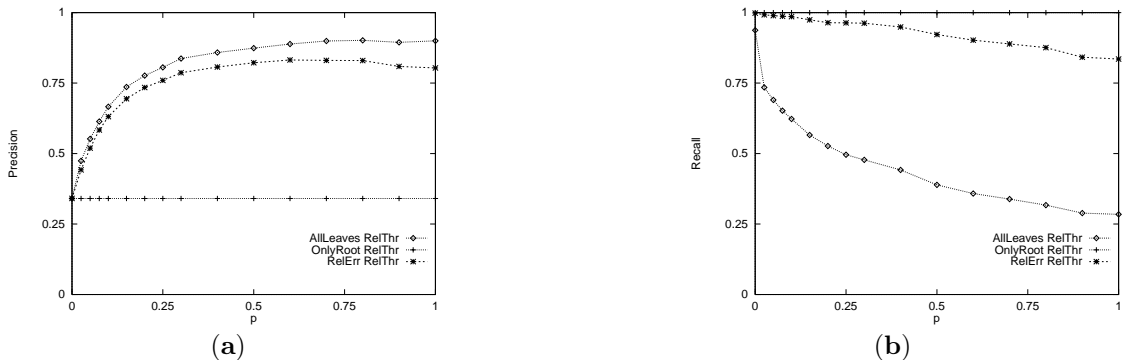


Figure 3: Average precision (a) and recall (b) for the link-based strategy of Section 3 as a function of  $p$  ( $\tau_c = 0.57$  for *RelErr*; *RelThr* pruning strategy).

for each of our *Spread* definitions and scope-pruning strategies. In Table 3, we report similar results for the content-based approach. In general, we see that the relative-threshold strategy to pruning scope works best in practice. In our data set, the content-based approach has a slight advantage over the link-based approach. However, we should regard these two approaches as complementary to each other for the following reasons, which we already touched on in Section 6.1. Often, web sites restrict robots from crawling their site (e.g., this is the case for The New York Times newspaper). In such cases, we cannot apply our content-based approach for estimating geographical scope, while we can still resort to the link-based approach. In other cases, the number of incoming links to a web site may be limited. In these cases, we should use our content-based approach, as long as any useful geographical information can be extracted from such resources that are not heavily cited.

## 7 A Geographically Aware Search Engine

Based on the techniques we developed in the previous sections, we have implemented a geographically aware search engine that downloads and indexes the full contents of 436 on-line newspapers based in the United States. Our search engine estimates the geographical

scope of the newspapers using the link-based technique of Section 3 with the *Entropy* definition for *Spread* and the *RelThr* scope-pruning strategy. This search engine is available at <http://www.cs.columbia.edu/~gravano/GeoSearch>.

Our search engine automatically pre-computes the geographical scope of the 436 newspapers that it indexes. When users query the engine, they specify their zip code in addition to their list of search keywords. Our system first uses just the keywords to *rank* the newspaper articles on those keywords using a standard, off-the-shelf text search engine called Swish. Our system then filters out all pages coming from newspapers whose geographical scope does not include the user’s specified zip code. Furthermore, our engine re-computes the score for each surviving page and returns the pages ranked in the resulting order. A page’s new score is a combination of the Swish-generated score for the page and the *Power* of the location in the geographical scope of the page’s newspaper that encloses the user’s zip code. Figure 5 shows the results for query “startups business” with zip code 94043, which corresponds to Mountain View, California. The first article is from The Nando Times, a national online newspaper. Our system has determined that this newspaper’s geographical scope is the whole country, hence the coloring of the map next to the correspond-

	<i>TopK</i>			<i>AbsThr</i>			<i>RelThr</i>		
	<i>F</i>	$\tau_c$	<i>k</i>	<i>F</i>	$\tau_c$	$\tau_e$	<i>F</i>	$\tau_c$	<i>p</i>
<i>AllLeaves</i>	0.34	–	$\infty$	0.51	–	0.0007	0.45	–	0.1
<i>OnlyRoot</i>	0.35	–	1	0.35	–	0	0.35	–	0
<i>VectorSpace</i>	0.76	0.7	1	0.52	0.9	0.0007	<b>0.81</b>	0.7	0.5
<i>Entropy</i>	0.78	0.8	1	0.51	0.8	0.0007	<b>0.82</b>	0.8	0.6
<i>RelErr</i>	0.78	0.57	1	0.52	0.67	0.0007	<b>0.83</b>	0.57	0.6

Table 2: Best average *F*-measure results for different *Spread* definitions (Section 3.2) and scope-pruning strategies (Section 3.3), using the link-based strategy of Section 3.

	<i>TopK</i>			<i>AbsThr</i>			<i>RelThr</i>		
	<i>F</i>	$\tau_c$	<i>k</i>	<i>F</i>	$\tau_c$	$\tau_e$	<i>F</i>	$\tau_c$	<i>p</i>
<i>AllLeaves</i>	0.37	–	1	0.42	–	0.0007	0.59	–	0.2
<i>OnlyRoot</i>	0.35	–	1	0.35	–	0	0.35	–	0
<i>VectorSpace</i>	<b>0.85</b>	0.6	1	0.82	0.6	0.1	<b>0.86</b>	0.6	0.7
<i>Entropy</i>	<b>0.85</b>	0.8	1	0.82	0.8	0.1	<b>0.85</b>	0.8	0.7
<i>RelErr</i>	0.76	0.57	1	0.72	0.50	0.1	0.76	0.57	0.5

Table 3: Best average *F*-measure results for different *Spread* definitions (Section 3.2) and scope-pruning strategies (Section 3.3), using the content-based strategy of Section 4.

ing article. The second article returned is from the San Jose Mercury News, a newspaper based in San Jose, California, whose technology reports have followers across the country. Our search engine has classified this newspaper as having a national geographical scope. The last article returned originated in a newspaper whose geographical scope consists of the entire state of California, which is marked with a solid color on the map, plus a few cities scattered across the country, indicated by placing a dot in their corresponding states.

## 8 Conclusion

In this paper, we discussed how to estimate the geographical scope of web resources, and how to exploit this information to build geographically aware applications. The main contributions of this paper include automatic estimation algorithms based on web-page content and HTML link information, metrics to evaluate the quality of such algorithms, a comprehensive evaluation of these techniques in a realistic experimental scenario, and an implementation of a geographically aware search engine for newspaper articles. One of the key observations of this paper is that the content-based techniques and the link-based techniques have specific advantages and disadvantages, and in fact can be used as complementary estimators of the scope of web resources. In effect, some sites might not allow us to “crawl” their contents, preventing us from using our content-based techniques. Some other sites might have a low number of incoming HTML links, preventing us from using our link-based techniques re-

liably. By combining these two approaches we can accurately estimate the geographical scope of many web resources, hence capturing a crucial dimension of web data that is currently ignored by search engines.

## Acknowledgments

This material is based upon work supported by the National Science Foundation under Grants No. IIS-97-33880 and IRI-96-19124. We also thank Jon Oringer for implementing the search engine of Section 7, and Jun Rao and Vasilis Vassalos for useful comments on the paper.

## References

- [1] S. Brin and L. Page. The anatomy of a large-scale hypertextual web search engine. In *Proceedings of the Seventh International World Wide Web Conference (WWW7)*, Apr. 1998.
- [2] O. Buyukkokten, J. Cho, H. García-Molina, L. Gravano, and N. Shivakumar. Exploiting geographical location information of web pages. In *Proceedings of the ACM SIGMOD Workshop on the Web and Databases (WebDB’99)*, June 1999.
- [3] J. P. Callan, Z. Lu, and W. B. Croft. Searching distributed collections with inference networks. In *Proceedings of the Eighteenth ACM International Conference on Research and Development in Information Retrieval (SIGIR’95)*, July 1995.



Figure 5: Search results from our geographically-aware search engine.

- [4] S. Chakrabarti, B. Dom, and P. Indyk. Enhanced hypertext categorization using hyperlinks. In *Proceedings of the 1998 ACM International Conference on Management of Data (SIGMOD'98)*, June 1998.
- [5] S. Chakrabarti, B. Dom, S. R. Kumar, P. Raghavan, S. Rajagopalan, A. Tomkins, D. Gibson, and J. Kleinberg. Mining the web's link structure. *IEEE Computer Magazine*, 32(8):60–67, 1999.
- [6] S. Chakrabarti, B. Dom, P. Raghavan, S. Rajagopalan, D. Gibson, and J. Kleinberg. Automatic resource compilation by analyzing hyperlink structure and associated text. In *Proceedings of the Seventh International World Wide Web Conference (WWW7)*, Apr. 1998.
- [7] D. Day, J. Aberdeen, L. Hirschman, R. Kozierok, P. Robinson, and M. Vilain. Mixed-initiative development of language processing systems. In *Proceedings of the Fifth ACL Conference on Applied Natural Language Processing*, Apr. 1997.
- [8] J. C. French, A. L. Powell, C. L. Viles, T. Emmitt, and K. J. Prey. Evaluating database selection techniques: A testbed and experiment. In *Proceedings of the Twentyfirst ACM International Conference on Research and Development in Information Retrieval (SIGIR'98)*, Aug. 1998.
- [9] D. Gibson, J. Kleinberg, and P. Raghavan. Inferring web communities from link topology. In *Proceedings of the Ninth ACM Conference on Hypertext and Hypermedia*, pages 225–234, June 1998.
- [10] L. Gravano, H. García-Molina, and A. Tomasic. *GLOSS: Text-source discovery over the Internet*. *ACM Transactions on Database Systems*, 24(2), June 1999.
- [11] R. W. Hamming. *Coding and Information Theory*. Prentice-Hall, 1980.
- [12] J. Kleinberg. Authoritative sources in a hyperlinked environment. In *Proceedings of the Ninth Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 668–677, Jan. 1998.
- [13] W. Meng, K.-L. Liu, C. T. Yu, X. Wang, Y. Chang, and N. Rishe. Determining text databases to search in the Internet. In *Proceedings of the Twenty-fourth International Conference on Very Large Databases (VLDB'98)*, Aug. 1998.
- [14] G. Salton. *Automatic Text Processing: The transformation, analysis, and retrieval of information by computer*. Addison-Wesley, 1989.
- [15] C. J. Van Rijsbergen. *Information Retrieval*. Butterworths, 1979.