

# Measuring Similarity to Detect Qualified Links

Xiaoguang Qi, Lan Nie, and Brian D. Davison

Dept. of Computer Science & Engineering  
Lehigh University



LEHIGH  
UNIVERSITY

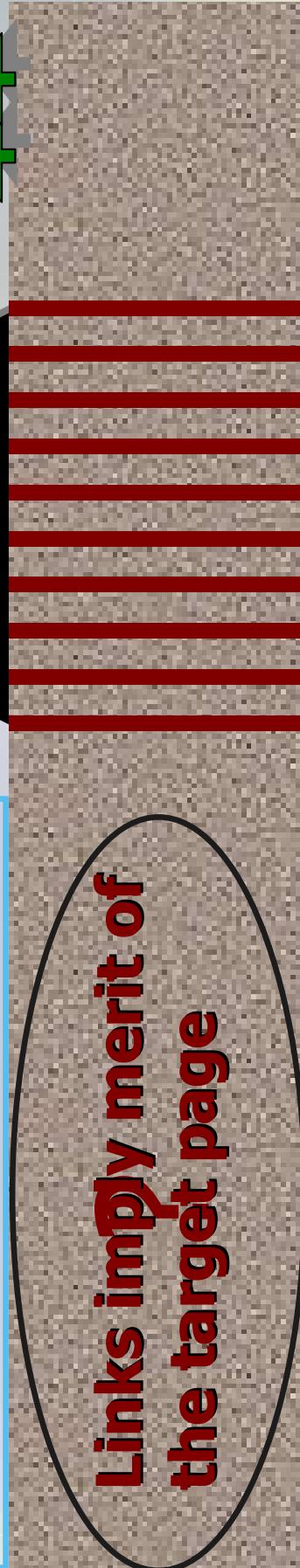
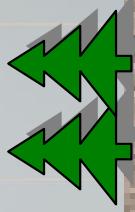
- **Introduction**
- **Approach**
- **Experiments**
- **Discussion & Conclusion**

Is this always true?  
**No!**

**Links imply merit of  
the target page**

Traditional link-based  
ranking algorithms

(PageRank, HITS, and  
most variations)

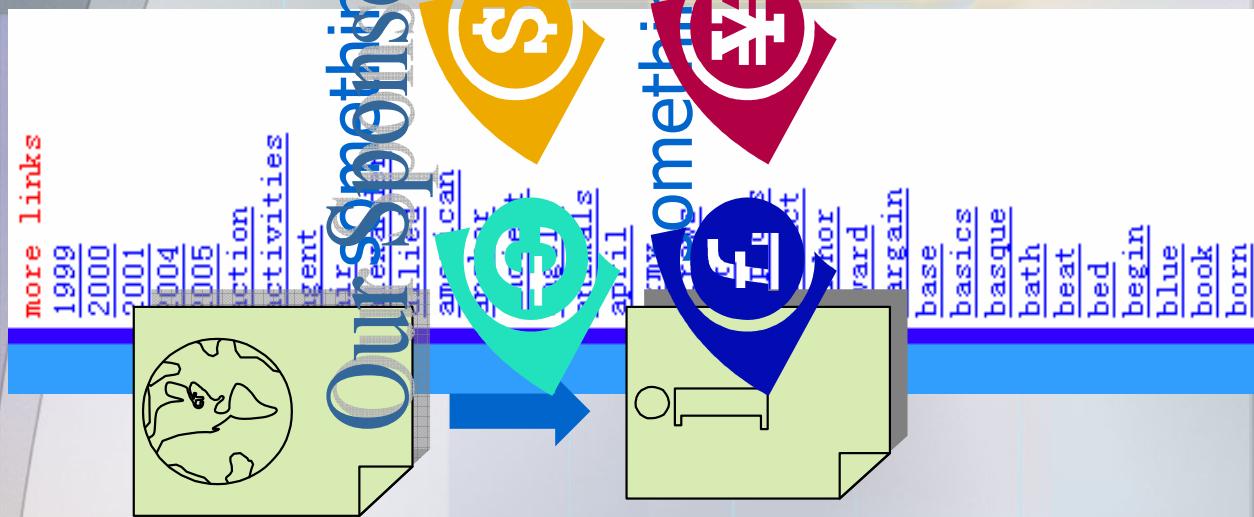


- **Spam links**

- **Navigational links**

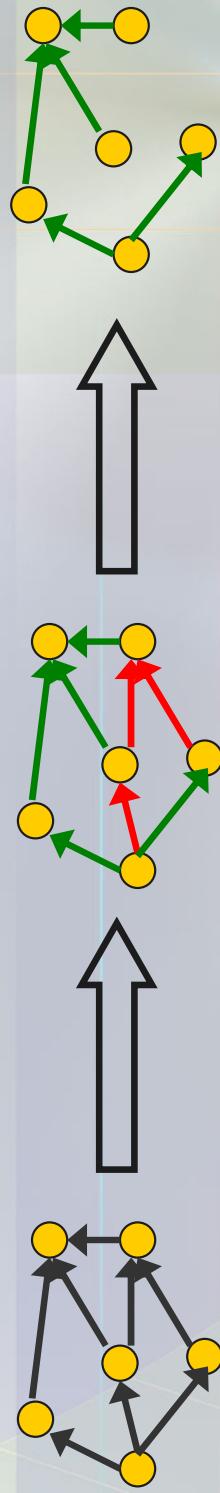
- **Advertising links**

- **Other irrelevant links**



- **These links**
  - may be useful for humans;
  - are effectively noise for link analysis.
- **Traditional link analysis algorithms do not distinguish them from useful links.**
  - As a result, the target pages of these links could get unmerited higher ranking.

- “Qualified links”
  - Links that are qualified to make a recommendation regarding the target page
- Our proposed approach
  1. Identify and filter out “unqualified links”
  2. Perform link analysis on the reduced link graph
- In our experiments, this approach can boost ranking performance.



# Background

- **Hyperlink-Induced Topic Search (HITS)**  
[Kleinberg 1997-1999]
  - The score of a hub (authority) depends on the sum of the connected authorities (hubs)
- **Bharat and Henzinger (1998)**
  - A number of improvements to HITS
    - imp*
  - Re-weight links involved in mutual reinforcement
  - Drop links within the same host

# Background (Cont.)

- **PageRank**

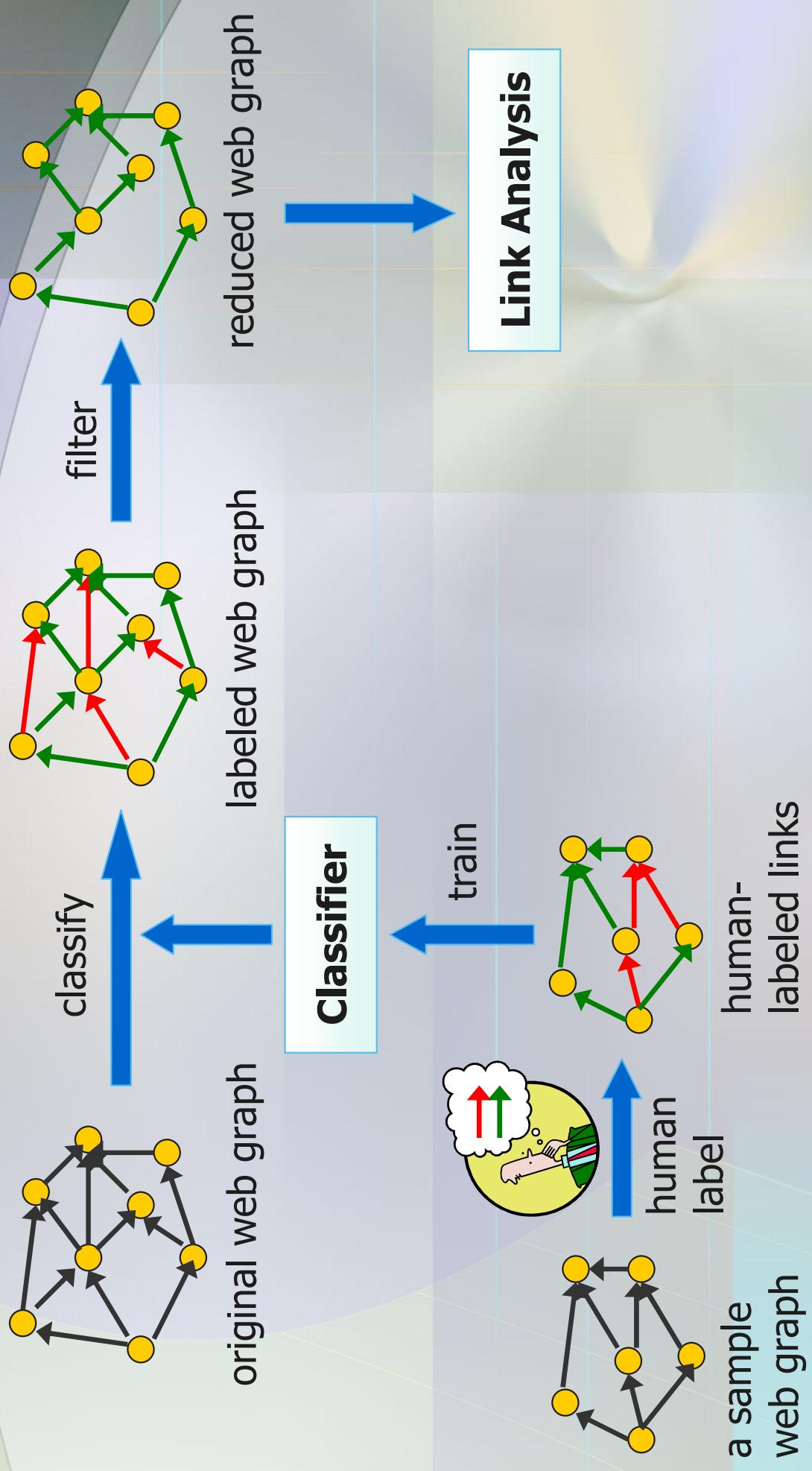
- **Random surfer model**

$$PR(i) = (1 - d) \sum_{j:j \rightarrow i} \frac{PR(j)}{O(j)} + d \frac{1}{N}$$

# Related Work

- Chakrabarti et al. (1998) improving HITS by adjusting the weights of links according to their anchor text and surrounding text.
- Davison (2000) proposed automatic recognition of nepotistic links using hand-crafted features.
- Lempel and Moran (2000) defined the tightly-knit community (TKC) effect.
- Li et al. (2002) pointed out the small-in, large-out communities could dominate HITS results.
- Wu and Davison (2005) proposed a two-step algorithms to identify link farms.
- Benczur et al. (2006) proposed to detect nepotistic links using language models.
- Carvallo et al. (2006) proposed to detect noisy links at the site level by examining link structure among web sites.

- Introduction
- Approach
- Experiments
- Discussion & Conclusion



# A navigational link

The screenshot shows the homepage of the ITT Pure-Flo website. At the top, there's a navigation bar with links for Home, Edit, View, File, and Members. Below this, there's a main menu with options like Valves, Careers, Contact, About Pure-Flo, Literature, News & Events, Members, and Search. A red circle highlights the ITT logo in the bottom left corner of the page. The main content area features a banner for the Fall 2006 Newsletter, a world map, and a section titled "Welcome To Pure-Flo". There are also "What's New" and "Find Sales & Service" buttons.

The screenshot shows the homepage of the ITT Corporation website. At the top, there's a navigation bar with links for Home, Company Profile, News, Investor Relations, Business & Products, Careers, and Contact Us. Below this, there's a search bar and a "GO" button. The main content area features a banner for "Advancing Human Progress" and another for "‘Green’ buildings in Beijing". On the right side, there's a "Latest News" section with several news items, a "ITT (NYSE)" section with stock market information, and a "Read our Stock Quote" section. A blue arrow points from the Pure-Flo logo on the left towards the ITT logo on the right.

This web site is designed, managed and published by Applegate Directory Ltd  
Copyright © 1996-2007 Applegate Directory Ltd  
Web Solutions by Abacus Tree  
Legal Privacy Contact Applegate

applegate

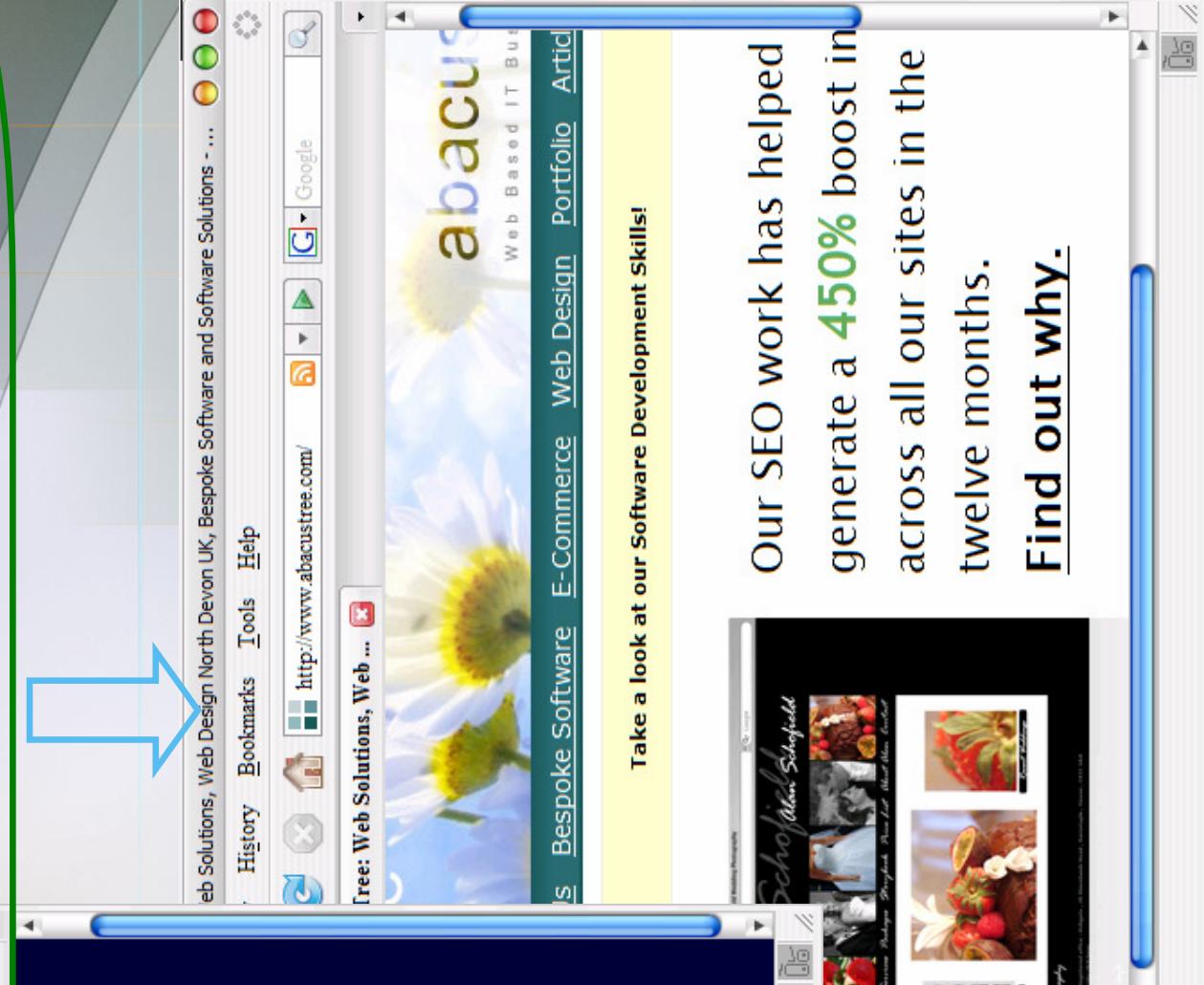
[www.applegate.co.uk](http://www.applegate.co.uk)

## Welcome to The Applegate Directory

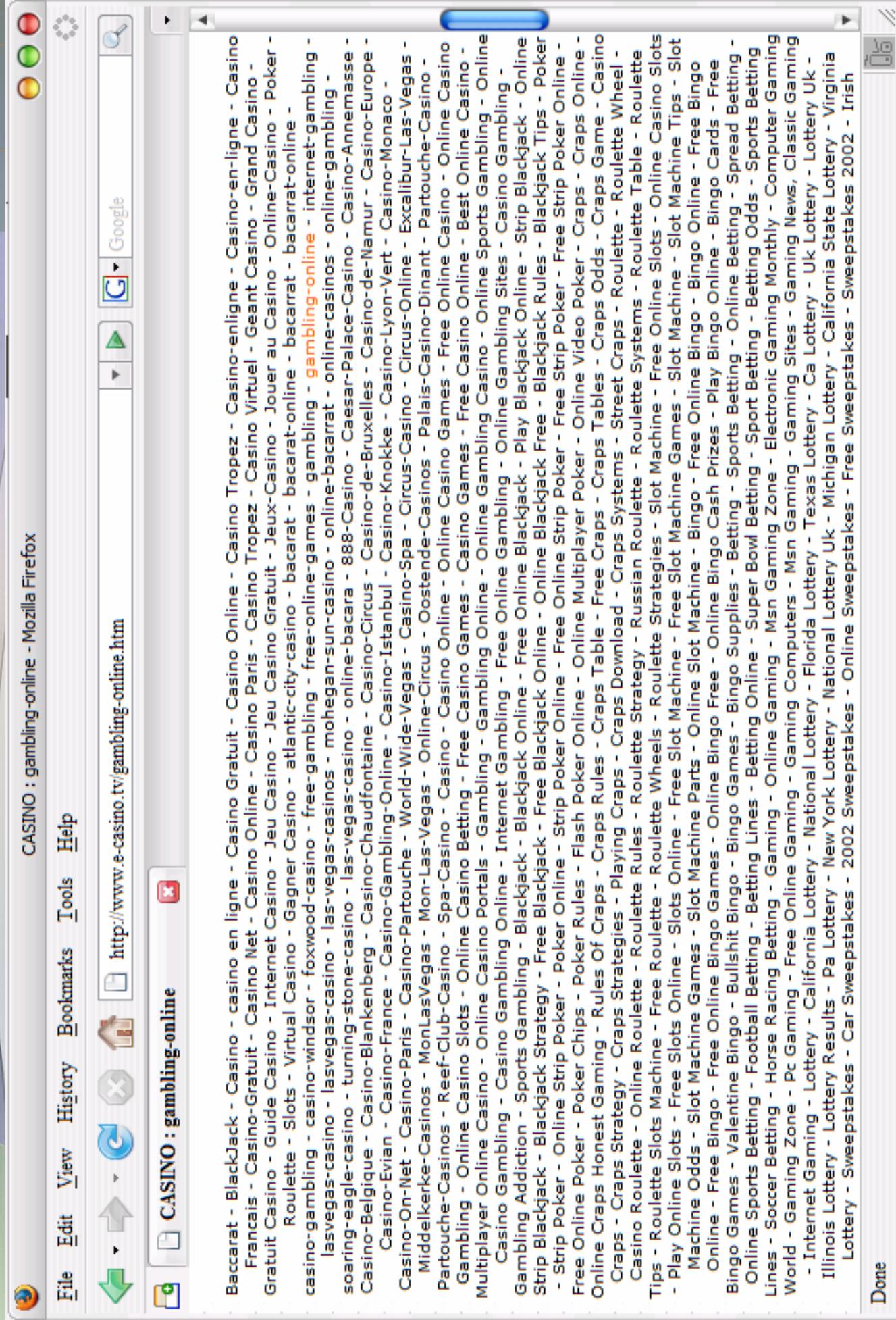
If you have come here from any domain other than  
[www.applegate.co.uk](http://www.applegate.co.uk) and wish to visit the Applegate Directory  
then please [click here](#).

This web site is designed, managed and published by Applegate Directory Ltd  
Copyright © 1996-2007 Applegate Directory Ltd  
Web Solutions by Abacus Tree  
Legal Privacy Contact Applegate

Done



d



# How to determine whether a link is qualified?

- There are many features that can be used.
- In this preliminary work, we studied six similarity measures between the source and target pages.
  - Hostname, URL, topic vector, tfidf content, anchor text, non-anchor text

# Similarity measures

- **Hostname similarity**

- The portion of common substrings of two hostnames

$$\text{Sim}_{\text{host}}(x, y) = \frac{2 \times |\text{Substr}(\text{host}_x, r) \cap \text{Substr}(\text{host}_y, r)|}{|\text{Substr}(\text{host}_x, r)| + |\text{Substr}(\text{host}_y, r)|}$$

- **URL similarity**

- Analogous to hostname similarity
- The portion of common substrings of two URLs

# Similarity measures (Cont.)

## •Topic vector similarity

- Cosine similarity of two topic vectors

$$\text{Sim}_{\text{topic}}(x, y) = \sum_{i=1}^n v_{xi} \times v_{yi}$$

-A topic vector  $v_x = (v_{x,1}, v_{x,2}, \dots, v_{x,n})$

- Each component is the probability that page  $x$  is on topic  $t$
- Can be computed using a classifier

## •Tfidf content similarity

- Cosine similarity of the two tfidf vectors

$$\sum_{t \in T} (x_t \times y_t)$$

$$\text{Sim}_{\text{content}}(x, y) = \frac{\sum_{t \in T} x_t^2 \times \sum_{t \in T} y_t^2}{\sqrt{\sum_{t \in T} x_t^2} \times \sqrt{\sum_{t \in T} y_t^2}}$$

# Similarity measures (Cont.)

- **Anchor text similarity**

- Analogous to content similarity
- Similarity (tfidf cosine) of anchor text on the two pages

- **Non-anchor text similarity**

- Analogous to content similarity
- Similarity of non-anchor text on the two pages

- Introduction
- Approach
- Experiments
- Discussion & Conclusion

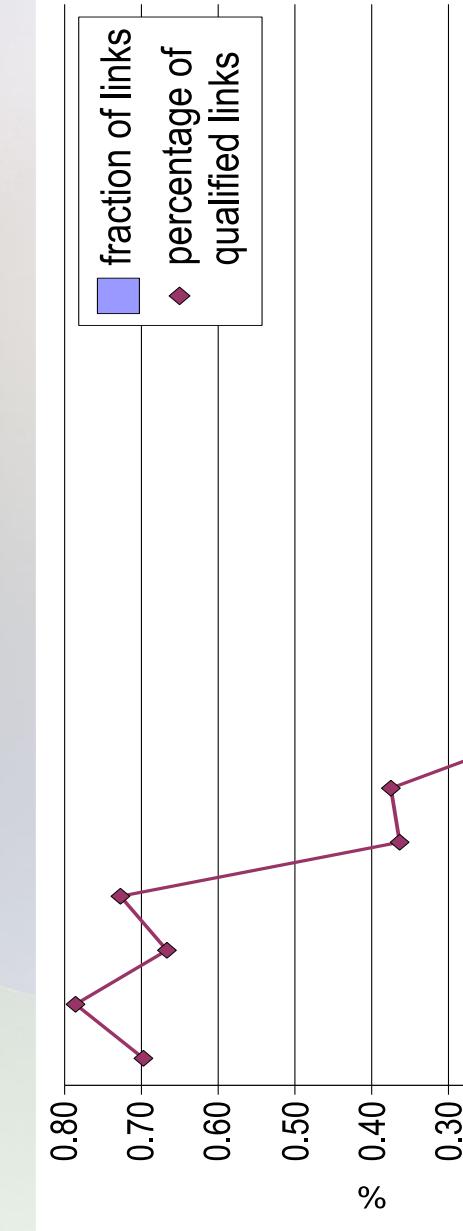
# Datasets

- **Query-specific datasets**
  - Collected and used by [Wu and Davison 2005]
  - 58 queries used in previous research, from ODP category names, and popular queries
- **Global dataset**
  - A 2005 crawl from the Stanford WebBase

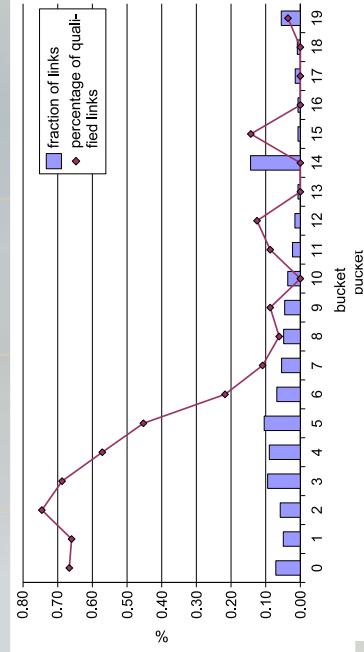
# Link classification

- **Human labeling of links**
  - More than a thousand links randomly selected from five query-specific datasets
  - Two human editors
- **Link classification based on all features**
  - A linear SVM classifier using SVM<sup>light</sup>
  - Two fold cross validation
  - Average accuracy 83.8%

# Link classification (Cont.)



Content similarity

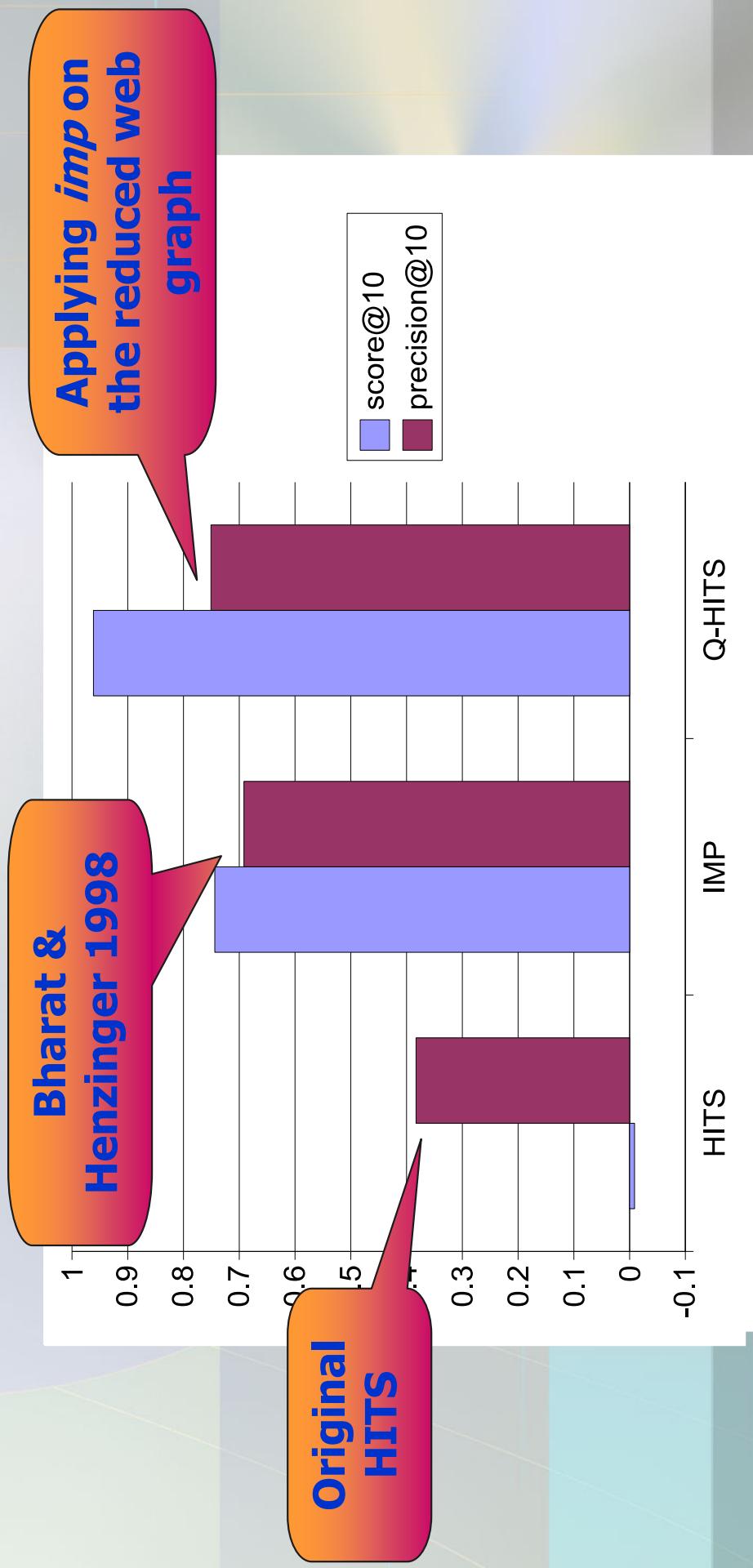


Non-anchor text similarity

Anchor text similarity was found to be the most discriminative feature.

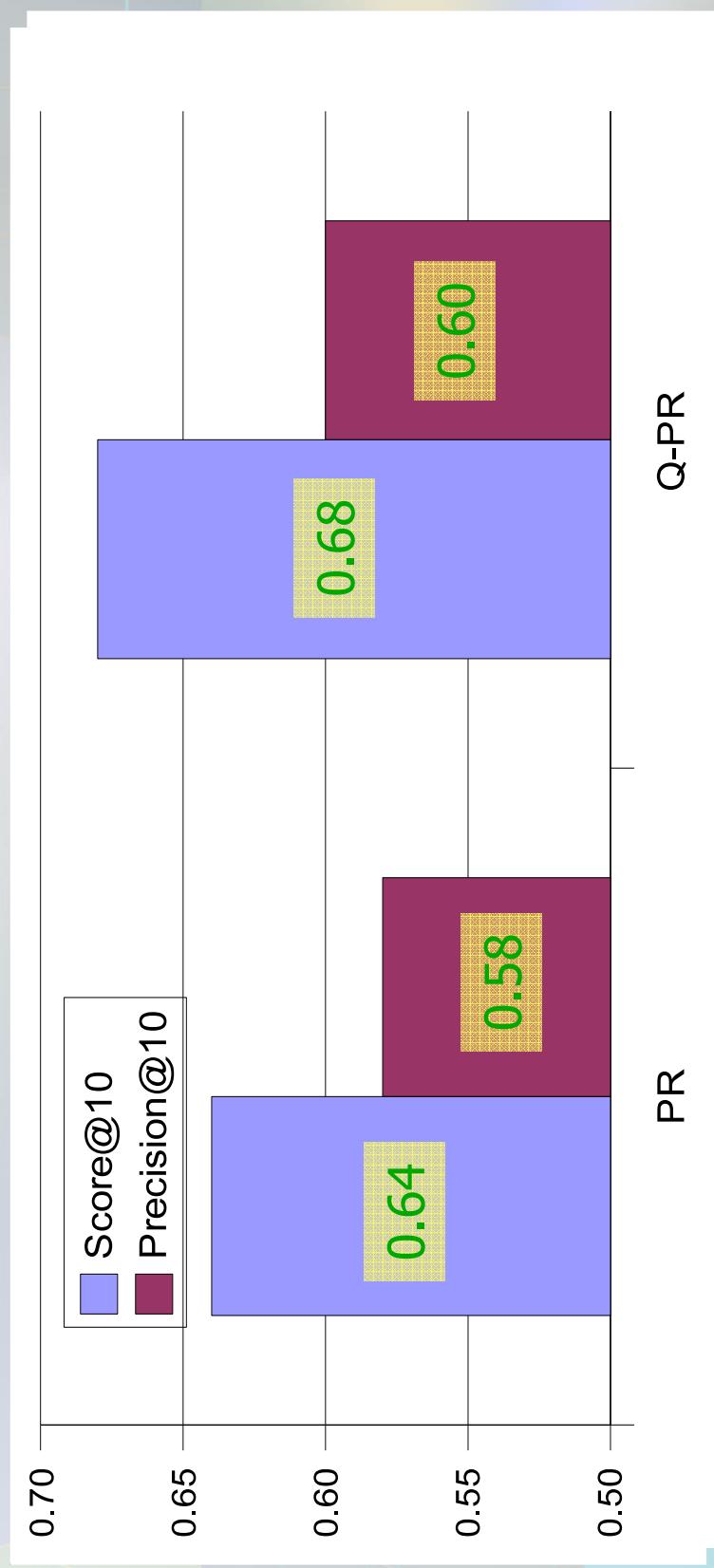
# Qualified HITS

- Performed on query-specific datasets.
- Unqualified links, identified by the classifier, are removed.
- *imp* is performed on the reduced link graph.
- In this experiment, 37% of the links are removed.



# Qualified PageRank

- Performed on WebBase dataset.
- Unqualified links, identified by the classifier, are removed.
- *PageRank* is performed on the reduced link graph.
- In this experiment, 0.4% of the links are removed.



- **Introduction**
- **Approach**
- **Experiments**
- **Discussion & Conclusion**

# Conclusion

- **Summary of approach**
  - Remove noise (unqualified links) from link graph before link analysis is applied
  - Identify unqualified links by measuring similarities between their source and target pages
- **Qualified HTTS is able to improve precision by 9% over Bharat and Henzinger's imp.**

# DISCUSSION

- How to determine the qualification of a link is an open question.
  - We used similarities between source and target page to demonstrate the potential of the idea.
  - What about other similarity measures?
  - What about non-similarity features?
- A closer look at link classification.
  - We trained a multi-class classifier to distinguish between qualified, spam, navigational, advertising, and other irrelevant links.
  - The classifier is effective in finding spam links
    - Not very helpful for other types.

# Thank you!



**WUIME Laboratory**  
<http://wuime.cse.lehigh.edu/>

**LEHIGH**  
UNIVERSITY

