



Tutorial:
Web Information Retrieval

Monika Henzinger

<http://www.henzinger.com/~monika>



What is this talk about?

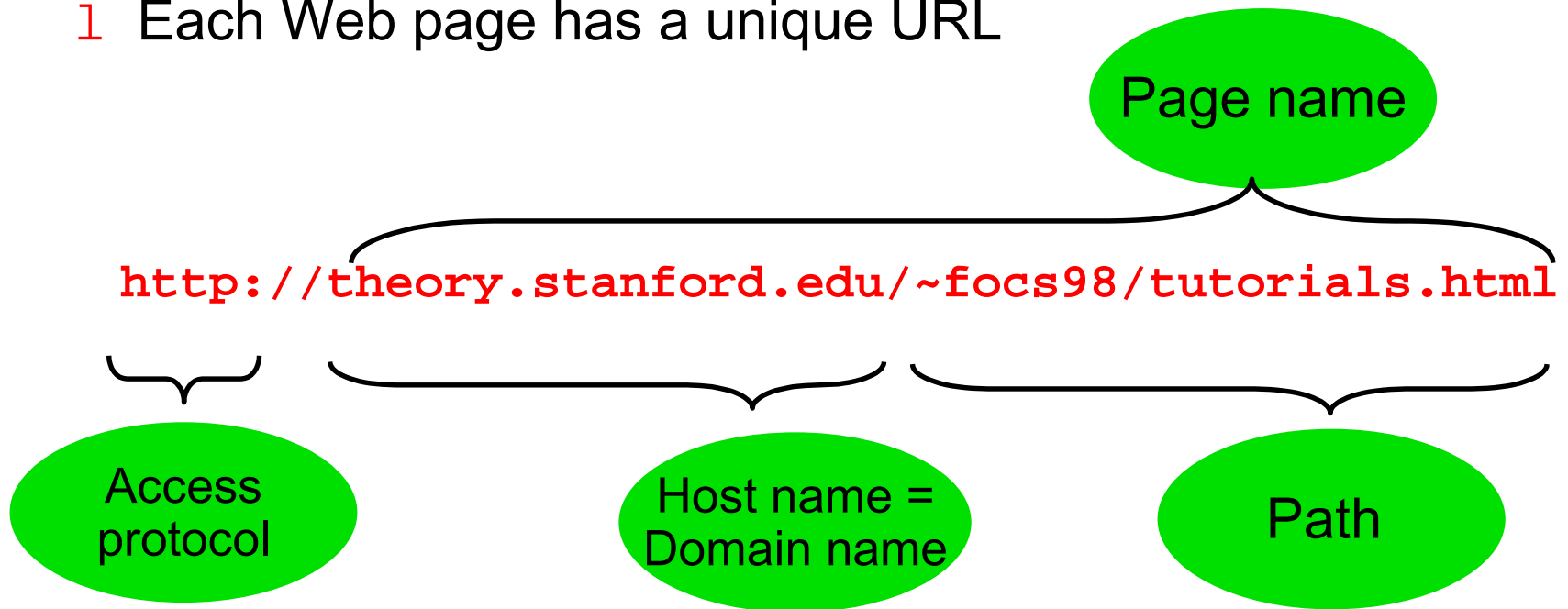
1 Topic:

- Algorithms for retrieving information on the Web

1 Non-Topics:

- Algorithmic issues in classic information retrieval (IR), e.g. stemming
- String algorithms, e.g. approximate matching
- Other algorithmic issues related to the Web:
 - networking & routing
 - cacheing
 - security
 - e-commerce

1 Each Web page has a unique URL



1 (Hyper) link = pointer from one page to another, loads second page if clicked on

1 In this talk: document = (Web) page

princess diana

Engine 1

[Princess Diana Memorial WebRing](#)
Follow the WebRing for a tour of memorial site
87% <http://www.geocities.com/RainForest/Vines/1009/diana1998>
[Grouped results from http://www.geocities.com](#)

[FOR DIANA, PRINCESS OF HEART - Dr. K](#)
...
Dr. Kate Wachs Comments on Princess Diana T
84% <http://www.therelationshipcenter.com/diana.shtml> (S)

[Princess Diana Editorial Cartoons! Cartoons :](#)
The Professional Cartoonists Index is the most c
cartoonists o
daily cartoon
82% <http://www>

[Diana, Princess of Wales](#)
1 July 1961 - 31 August 1997 The BBC Web sit
Camera Press/Snowdon
79% <http://www.royal.gov.uk/start.htm> (Size 2.3K) Doc
[Grouped results from http://www.royal.gov.uk](#)

Relevant and high quality

Engine 2

1. [Re: Lost in the shadow of Princess Diana](#)
[URL: www.spiceisle.com/talkshop/messages/6232.htm]
The SpiceIslander TalkShop. [Follow Ups] [Pos
The SpiceIslander TalkShop] Date: September
00:54:03 From: Sno,...
Last modified 12-Sep-97 - page size 4K - in English [Tran

2. [Re: Princess Diana's gown auction](#)
[URL: www.elle.com/textes/blablaba/forum/messages1/14
Re: Princess Diana's gown auction. [Follow Ups
Followup] [Elle International - Blablaba] Posted
September 07, 1997 at 02:15:26:..
Last modified 30-Mar-98 - page size 2K - in English [Tran

3. [Re: Princess Diana](#)
[URL: spicyhot.com/gaynet/messages/1053.html]
Re: Prince
Maine Ga
Novembe
Last modif

4. [Re: Princess Diana - Queen of Hearts](#)
[URL: www.elle.com/textes/blablaba/forum/messages1/26
Re: Princess Diana - Queen of Hearts. [Follow U
Followup] [Elle International - Blablaba] Posted
on August 31, 1997 at...
Last modified 30-Mar-98 - page size 4K - in English [Tran

Relevant but low quality

Engine 3

1. [Free Passwords To Adult Sites ...](#)
99% - **Articles & General info:** Free Passwords
Sites warez princess diana demi moore
magazine kathy ireland lingerie jennifer aniston cook
warez princess diana demi moore... 03/09/98
Commercial site: <http://www.purient.com/warez>

2. [SEX CHAT XXX NUDE PORN PLAYBOY P](#)
[AMERICAN FIRST FREE PICTURES WOMEN](#)
99% - **Articles & General info:** SEX CHAT X
PORN PLAYBOY PAMELA ANDERSON P
PICTURES WOMEN ADULT MUSIC CHAT P
BRITICA BERRY MCCARTHY LIPSONE SA
CHAT CRAWFORD STEE GIBLL... 03/09/98

Personal page: <http://www.connix.com/~wgonzo/sex/slidesuperall.htm>

3. [Ro](#)

Personal page: <http://www.octet.com/~gonzo/jy>

4. [Sunday, 18-Jan-98](#)
99% - **Articles & General info:** Sunday, 18-Jan-
CHAT XXX NUDE PORN PLAYBOY PAME

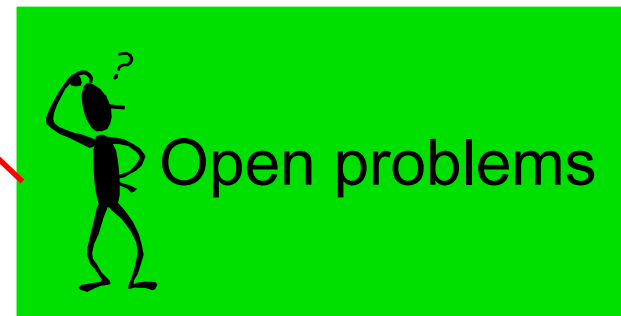
Not relevant index pollution

Google™ Outline

- 1 Classic IR vs. Web IR
- 1 Some IR tools specific to the Web
 - For each type
 - Examples
 - Algorithmic issues
- 1 Conclusions

Details on

- Ranking
- Duplicate elimination
- Search-by-example
- Measuring search engine index quality



- 1 **Input:** Document collection
- 1 **Goal:** Retrieve documents or text with information content that is **relevant** to user's **information need**
- 1 **Two aspects:**
 - 1. Processing the collection
 - 2. Processing queries (searching)
- 1 Reference Texts: SW'97, BR'99

“model” = strategy for determining which documents to return

- 1 Logical model: String matches plus AND, OR, NOT
- 1 Vector model (Salton et al.):
 - Documents and query represented as vector of terms
 - Vector entry i = weight of term i = function of frequencies within document and within collection
 - Similarity of document & query = cosine of angle of their vectors
 - Query result: documents ordered by similarity
- 1 Other models used in IR but not discussed here:
 - Probabilistic model, cognitive model, ...

- 1 **Input:** The publicly accessible Web
- 1 **Goal:** Retrieve **high quality** pages that are **relevant** to user's **need**
 - Static (files: text, audio, ...)
 - Dynamically generated on request: mostly data base access
- 1 **Two aspects:**
 1. Processing and representing the collection
 - Gathering the static pages
 - “Learning” about the dynamic pages
 2. Processing queries (searching)

(1) Pages:

- 1 Bulk >1B (12/99) [GL'99]
- 1 Lack of stability..... Estimates: 23%/day, 38%/week [CG'99]
- 1 Heterogeneity
 - Type of documents .. Text, pictures, audio, scripts,...
 - Quality From dreck to ICDE papers ...
 - Language 100+
- 1 Duplication
 - Syntactic..... 30% (near) duplicates
 - Semantic..... ??
- 1 Non-running text..... many home pages, bookmarks, ...
- 1 High linkage..... ≥ 8 links/page in the average

Typical home page: non-running text

COMPAQ

Q
Better answers.sm

find out more

NEWS

- [Compaq AlphaServer Awarded Best Data Warehousing Server](#)
- [Celera, Human Genomic Research Company, Selects Compaq as IT Solutions, Services Partner](#)
- [Greater Performance from Compaq's New Generation of 64-bit AlphaServers](#)
- [Compaq Reports Third Quarter Results](#)

October 26, 1998

SEARCH | SITE GUIDE | COMMENTS

BUY ONLINE
buy direct from us & other ways to shop

PRODUCTS
products and solutions for home & business computing

SERVICES
programs, support & downloadables

BUSINESS COMPUTING
products training & solutions for companies

HOME COMPUTING
products, services & online shopping for home computing

RESELLERS & PARTNERS
information for our ASEs, sales & business partners

GOVERNMENT, EDUCATION & HEALTHCARE
for our customers in federal, state and local government, K-12 and higher education, and healthcare

INFO MESSENGER
stay on track with product updates

INSIDE COMPAQ
corporate information, news, publications, jobs and key contacts at Compaq

YEAR 2000
Compaq products & year 2000 compliance information

DIGITAL.COM | TANDEM.COM

Pick a country site

LEGAL NOTICES AND PRIVACY STATEMENT



Typical home page: Non-running text



Welcome to **Embassy of India** Washington DC

[Archive](#) | [Calendar](#) | [Gallery](#) | [Links](#) | [Site Index](#)

Embassy

Press Releases

Consular Services

Publications

India Info

Foreign Relations

Policy Statements

India - US Relations

South Asia Region

What's New

[President's address to Parliament](#) - October 25, 1999

[Address to the nation by Prime Minister Atal Bihari Vajpayee](#)
[Vajpayee - A Profile](#) | [Council of Ministers](#)

[Cross Border Terrorism](#)

[Indian Parliament Elections - 1999](#) | [Results](#)

[Draft Report on Nuclear Doctrine](#) | [Daily News from Indian Media](#)

[Current Issues on Jammu & Kashmir](#) | [Kargil Situation](#)

[Join the Embassy's elist & discussion group](#)

This site is created and maintained by the Press & Information Wing,
Embassy of India. Comments on the website to webmaster@indiaqov.org



The big challenge

**Meet the user needs
given
the heterogeneity of Web pages**

(2) Users:

1 Make poor queries

- Short (2.35 terms avg)
- Imprecise terms
- Sub-optimal syntax (80% queries without operator)
- Low effort

1 Wide variance in

- Needs
- Knowledge
- Bandwidth

1 Specific behavior

- 85% look over one result screen only
- 78% of queries are not modified
- Follow links
- See various user studies in CHI, Hypertext, SIGIR, etc.



The bigger challenge

**Meet the user needs
given
the heterogeneity of Web pages
and
the poorly made queries.**

Why don't the users get what they want?

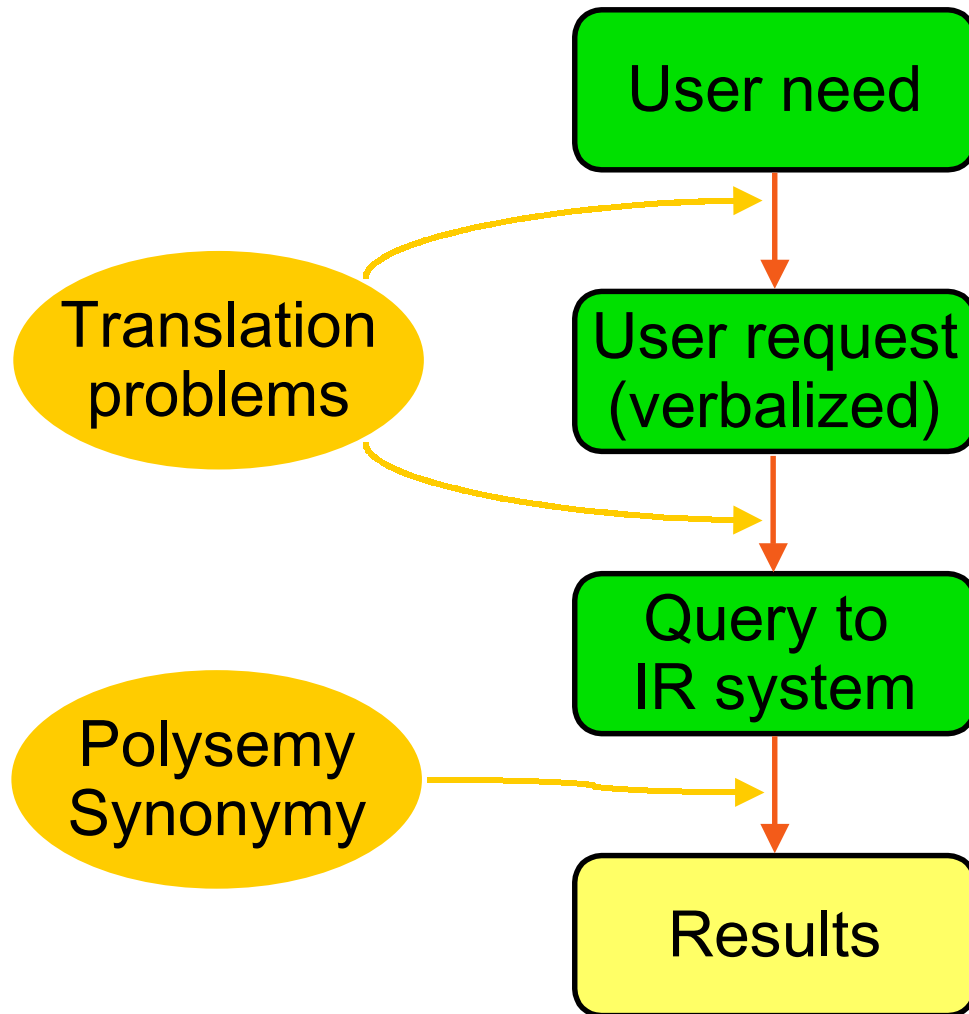
Example

I need to get rid of mice in the basement

What's the best way to trap mice alive?

`mouse trap`

Software, toy cars, inventive products, etc



bigsun.wbs.net/homepages/m/o/u/mouse_trap

New! Try out [GoogleScout](#)

[Doc Fizzix - Mousetrap Cars and Mouse Trap Powered Vehicles](#)

...The best mousetrap cars& **mouse trap** cars site! **Mouse...**

www.docfizzix.com/ [Cached \(14k\)](#) New! Try out [GoogleScout](#)

[Mouse Trap](#)

... **Mouse Trap Mouse Trap** is a simple but effective...

...can also be configured to **trap** the **mouse** on system startup or at a...

www.homeonthewww.com/ryan/mousetrap.html [Cached \(5k\)](#) New! Try out [GoogleScout](#)

[Tin Cat Repeating Mouse Trap](#)

www.biconet.com/critter/tincat.html [Cached \(11k\)](#) New! Try out [GoogleScout](#)

J

Smart Mouse Trap

www.biconet.com/critter/smt.html [Cached \(11k\)](#) [New!](#) Try out [GoogleScout](#)

J

Tin Cat Repeating Mouse Trap

www.biconet.com/critter/tincat.html [Cached \(11k\)](#) [New!](#) Try out [GoogleScout](#)

J

Horned Owl Inflatable Scarecrow

www.biconet.com/critter/owl.html [Cached \(10k\)](#) [New!](#) Try out [GoogleScout](#)

Rat Traps, mice traps ,glue boards, moth traps,pantry pest traps

...MULTIPLE TRAPS FOR **MICE** MOUSE MASTER A multiple catch **trap** for...

...Single **Trap** #855 + Lure \$22.06 SNAP TRAPS FOR RATS AND **MICE** RAT...

doyourownpestcontrol.com/traps.htm [Cached \(40k\)](#) [New!](#) Try out [GoogleScout](#)

J

Mice

...the wall (see our **trap** placement guide) since **mice** mostly navigate...

...those signs of **mice**? That's where you place the **trap**. **Mice**...

www.unexco.com/mice.html [Cached \(15k\)](#) [New!](#) Try out [GoogleScout](#)

J

National Food Safety Database: Disaster Handbook

...traps are needed in a house to **trap mice** than rats. Rats and...

...week before they approach a **trap**. **Mice** are curious and will...

www.foodsafety.org/dh/dh044.htm [Cached \(21k\)](#) [New!](#) Try out [GoogleScout](#)



The bright side: Web advantages vs. classic IR

User

- 1 Many tools available
- 1 Personalization
- 1 Interactivity (refine the query if needed)

Collection/tools

- 1 Redundancy
- 1 Hyperlinks
- 1 Statistics
 - Easy to gather
 - Large sample sizes
- 1 Interactivity (make the users explain what they want)

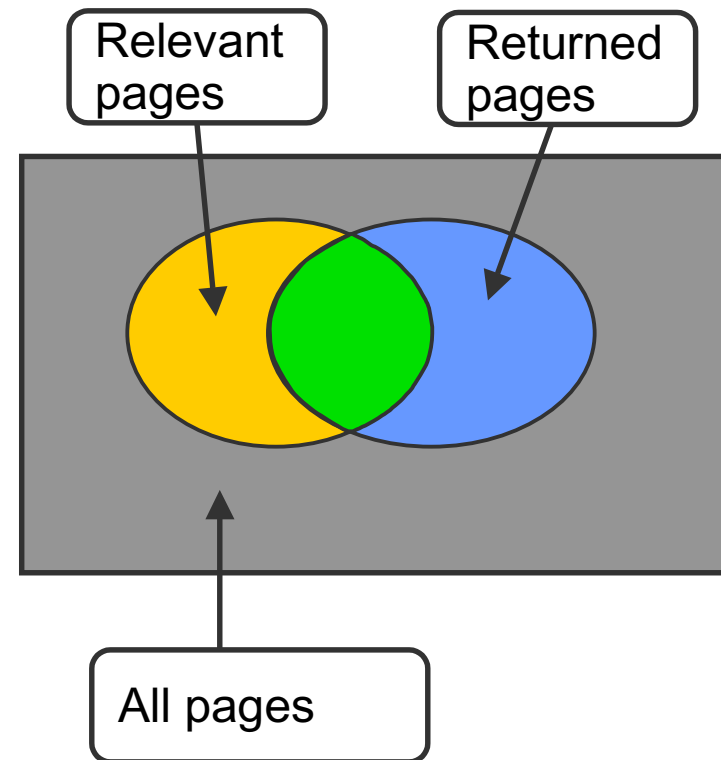
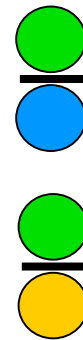


Quantifying the quality of results

- 1 How to evaluate different strategies?
- 1 How to compare different search engines?

We start from a **human made relevance judgement** for each **(query, page)** pair and compute:

- 1 **Precision**: % of returned pages that are relevant.
- 1 **Recall**: % of relevant pages that are returned.
- 1 **Precision at (rank) 10**: % of top 10 pages that are relevant
- 1 **Relative Recall**: % of relevant pages *found by some means* that are returned





Evaluation in the Web context

- 1 Quality of pages varies widely
- 1 We need both relevance and high quality = **value** of page.
- 1 **Precision at 10**: % of top 10 pages that are **valuable**
- ...

1 General-purpose search engines:

- direct: AltaVista, Excite, Google, Infoseek, Lycos,
- Indirect (Meta-search): MetaCrawler, DogPile, AskJeeves, InvisibleWeb, ...

1 Hierarchical directories: Yahoo!, all portals.

1 Specialized search engines:

- Home page finder: Ahoy
- Shopping robots: Jango, Junglee, ...
- Applet finders

Database
mostly built
by hand

Google™ *Web IR tools (cont...)*

- 1 **Search-by-example:** Alexa's "What's related", Excite's "More like this", Google's "Googlescout", etc.
 - 1 **Collaborative filtering:** Firefly, GAB, ...
 - 1 ...
-
- 1 **Meta-information:**
 - Search Engine Comparisons
 - Query log statistics
 - ...



General purpose search engines

1 Search engines' components:

- **Spider = Crawler** -- collects the documents
- **Indexer** -- process and represents the data
- **Search interface** -- answers queries



Algorithmic issues related to search engines

1 Collecting documents

- Priority
- Load balancing
 - Internal
 - External
- Trap avoidance
- ...

1 Processing and representing the data

- Query-independent ranking
- Graph representation
- Index building
- Duplicate elimination
- Categorization
- ...

1 Processing queries

- Query-dependent ranking
- Duplicate elimination
- Query refinement
- Clustering
- ...

1 **Goal:** order the answers to a query in decreasing order of value

- **Query-independent:** assign an intrinsic value to a document, regardless of the actual query
- **Query-dependent:** value is determined only wrt a particular query.
- **Mixed:** combination of both valuations.

1 **Examples**

- **Query-independent:** length, vocabulary, publication data, number of citations (indegree), etc.
- **Query-dependent:** cosine measure



Some ranking criteria

- 1 **Content-based** techniques (variant of term vector model or probabilistic model) – mostly query-dependent
- 1 **Ad-hoc factors** (anti-porn heuristics, publication/location data, ...) – mostly query-independent
- 1 **Human annotations**
- 1 **Connectivity-based** techniques
 - Query-independent: PageRank [PBMW'98, BP'98], indegree [CK'97], ...
 - Query-dependent: HITS [K'98], ...

- 1 Idea: Mine hyperlink information of the Web
- 1 Assumptions:
 - Links often connect related pages
 - A link between pages is a recommendation
- 1 Classic IR work (citations = links) a.k.a. “Bibliometrics”
[K’63, G’72, S’73,...]
- 1 Socio-metrics [K’53, MMSM’86,...]
- 1 Many Web related papers build on this idea [PPR’96, AMM’97, S’97, CK’97, K’98, BP’98,...]



Graph representation for the Web

- 1 A **node** for each page u
- 1 A **directed edge** (u,v) if page u contains a hyperlink to page v .



Query-independent ranking: Motivation for PageRank

- 1 Assumption: A **link** from page A to page B is a **recommendation** of page B by the author of A (we say B is *successor* of A)
- ⌚ Quality of a page is related to its in-degree

- 1 Recursion: Quality of a page is related to
 - its in-degree, and to
 - the quality of pages linking to it
- ⌚ **PageRank** [BP '98]

- 1 Consider the following infinite **random walk** (surf):
 - Initially the surfer is at a random page
 - At each step, the surfer proceeds
 - to a randomly chosen web page with probability d
 - to a randomly chosen successor of the current page with probability $1-d$
- 1 **The PageRank of a page p is the fraction of steps the surfer spends at p in the limit.**

Google PageRank (cont.)

Said differently:

- 1 Transition probability matrix is

$$d \times U + (1 - d) \times A$$

where U is the uniform distribution and A is adjacency matrix (normalized)

- 1 PageRank = stationary probability for this Markov chain, i.e.

$$PageRank(u) = \frac{d}{n} + (1 - d) \sum_{(v,u) \in E} PageRank(v) / outdegree(v)$$

where n is the total number of nodes in the graph

- 1 Used as one of the ranking criteria in Google



Output from Google: princess diana

[Diana, Princess of Wales](#)

... **Diana, Princess** of Wales 1 July 1961 - 31 August 1997 The BBC Web...
www.royal.gov.uk/start.htm [Cached \(2k\)](#) [New!](#) Try out [GoogleScout](#)

J

www.royal.gov.uk/

[New!](#) Try out [GoogleScout](#)

J

[Princess Diana: Remember Diana, Princess of Wales](#)

...This ribbon is in memory of **Diana, Princess** of Wales. Please put...
...indefinitely as a tribute to **Diana, Princess** of Wales. I feel I...
www.gargaro.com/diana.html [Cached \(8k\)](#) [New!](#) Try out [GoogleScout](#)

J

www.geocities.com/RainForest/Vines/1009/diana.htm

[New!](#) Try out [GoogleScout](#)

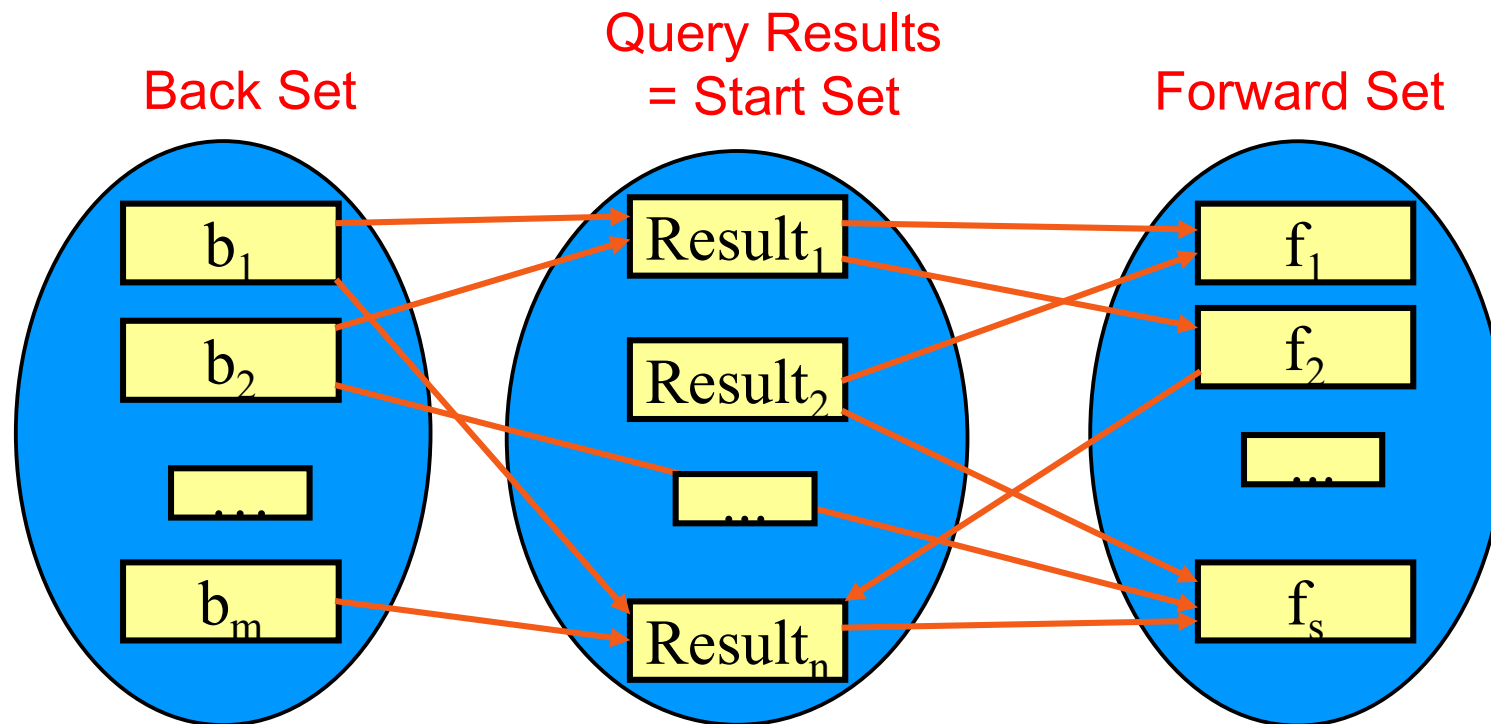
J

[CNN - The Death of Princess Diana](#)

...▪ The Burial: **Princess Diana's** coffin is taken to family...
...service held in memory of **Princess Diana** - VXtreme streaming video...
www.cnn.com/WORLD/9708/31/diana.links/ [Cached \(12k\)](#) [New!](#) Try out [GoogleScout](#)

Query-dependent ranking: the neighborhood graph

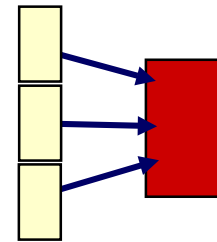
- 1 Subgraph associated to each query



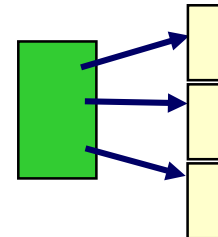
An edge for each hyperlink, but no edges within the same host

1 **Goal:** Given a query find:

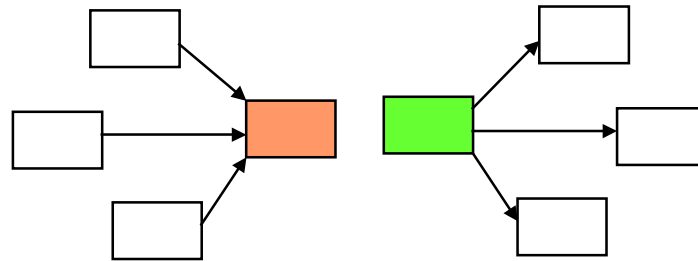
– Good sources of content (authorities)



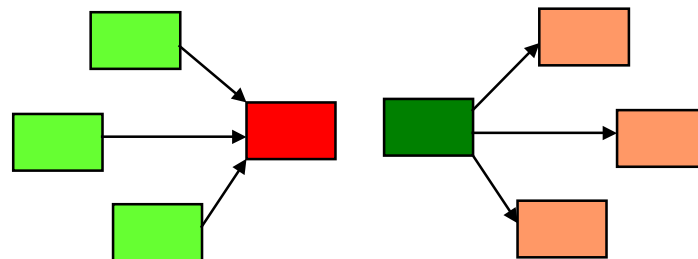
– Good sources of links (hubs)



- 1 **Authority** comes from in-edges.
Being a **good hub** comes from out-edges.



- 1 **Better authority** comes from in-edges from **good hubs**.
Being a **better hub** comes from out-edges to **good authorities**.

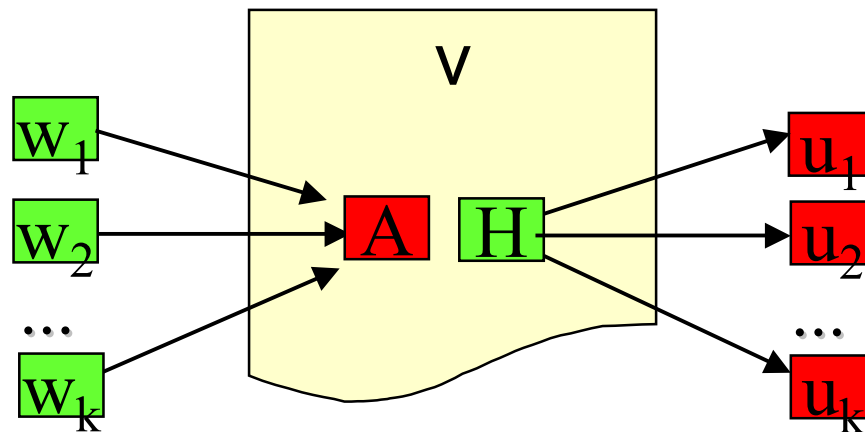


Repeat until \overrightarrow{HUB} and \overrightarrow{AUTH} converge:

Normalize \overrightarrow{HUB} and \overrightarrow{AUTH}

$HUB[v] := \sum AUTH[u_i]$ for all u_i with $Edge(v, u_i)$

$AUTH[v] := \sum HUB[w_i]$ for all w_i with $Edge(w_i, v)$





Output from HITS: jobs

1. www.ajb.dni.uk - British career website J
2. www.britnet.co.uk/jobs.htm
3. www.monster.com - US career website J
4. www.careermosaic.com - US career website J
5. plasma-gate.weizmann.ac.il/Job...
6. www.jobtrak.com - US career website J
7. www.occ.com - US career website J
8. www.jobserve.com - US career website J
9. www.allny.com/jobs.html - jobs in NYC
10. www.commart.com/bin/... - US career website J



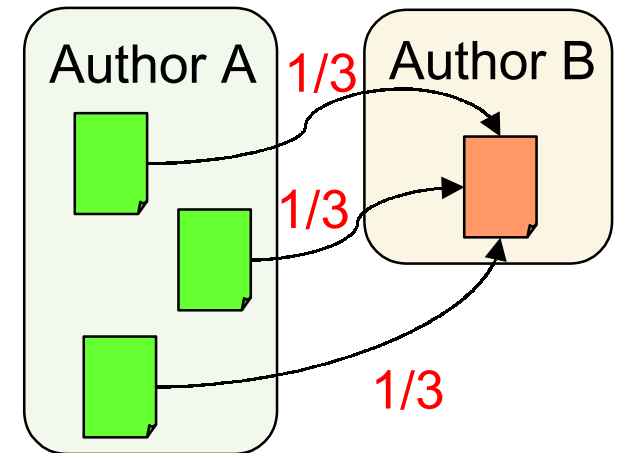
Output from HITS: +jaguar +car

1. www.toyota.com
2. www.chryslercars.com
3. www.vw.com
4. www.jaguravehicles.com J
5. www.dodge.com
6. www.usa.mercedes-benz.com
7. www.buick.com
8. www.acura.com
9. www.bmw.com
10. www.honda.com

Google™ *Problems & solutions*

- 1 Some edges are “wrong” -- not a recommendation:
 - multiple edges from same author
 - automatically generated
 - spam, etc.

Solution: Weight edges to limit influence



- 1 Topic drift
 - Query: **+jaguar +car**
 - Result: pages about **cars** in general

Solution: Analyze content and assign topic scores to nodes

Repeat until \vec{HUB} and \vec{AUTH} converge:

Normalize \vec{HUB} and \vec{AUTH}

$$HUB[v] := \sum AUTH[u_i] TopicScore[u_i] weight[v, u_i]$$

for all u_i with Edge(v, u_i)

$$AUTH[v] := \sum HUB[w_i] TopicScore[w_i] weight[w_i, v]$$

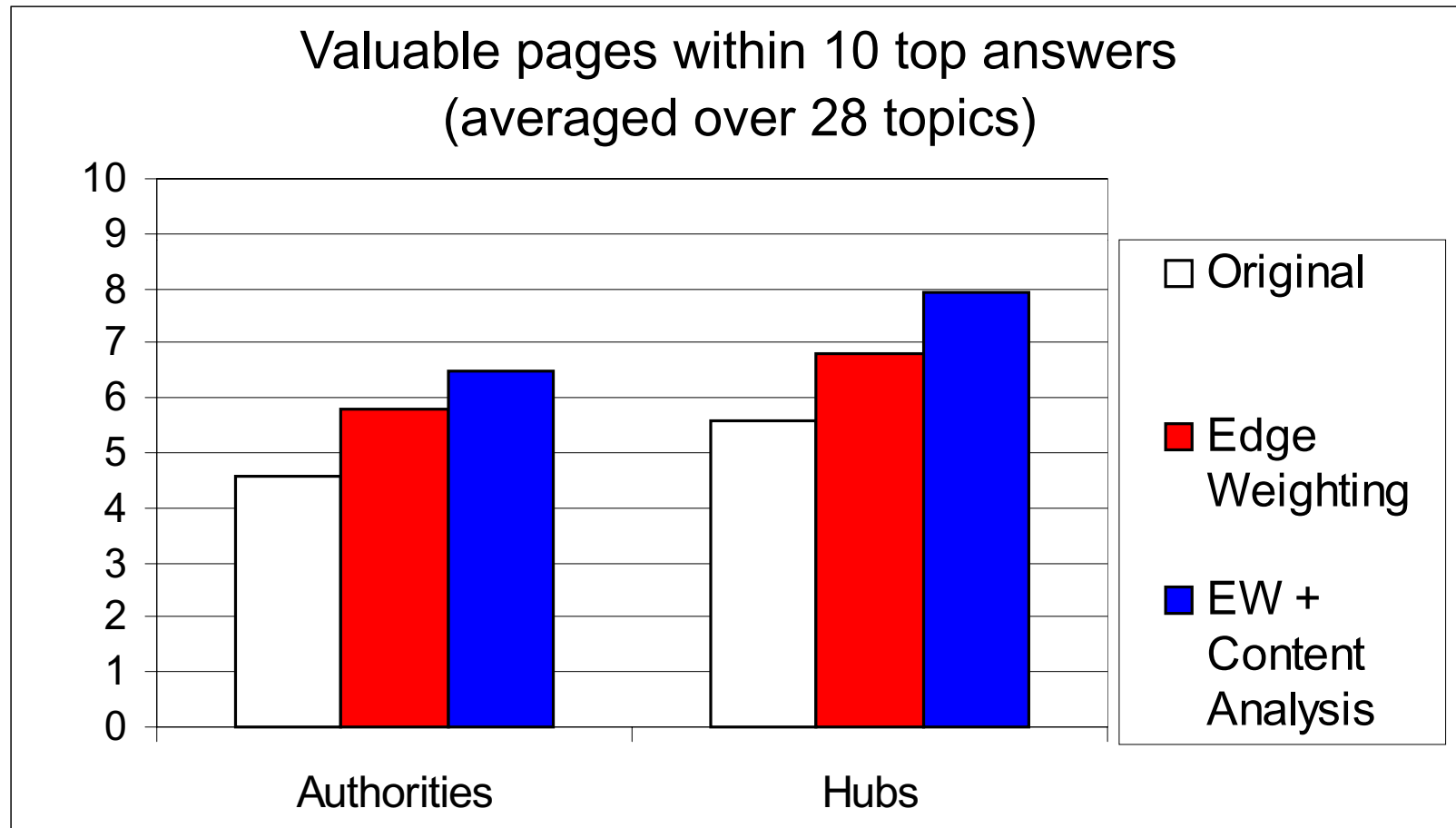
for all w_i with Edge(w_i, v)

[CDRRGK'98, BH'98, CDGKRRT'98]



Output from modified HITS: **+jaguar +car**

1. www.jaguarcars.com/ - official website of Jaguar cars J
2. www.collection.co.uk/ - official Jaguar accessories J
3. home.sn.no/.../jaguar.html - the Jaguar Enthusiast Place J
4. www.terrystjag.com/ - Jaguar Parts J
5. www.jaguarvehicles.com/ - official website of Jaguar cars J
6. www.jagweb.com/ - for companies specializing in Jags. J
7. jagweb.com/jdht/jdht.html - articles about Jaguars and Daimler
8. www.jags.org/ - Oldest Jaguar Club J
9. connection.se/jagsport/ - Sports car version of Jaguar MK II
10. users.aol.com/.../jane.htm -Jaguar Association of New England Ltd.



- | | |
|--|---|
| <ul style="list-style-type: none">1 Computation:<ul style="list-style-type: none">– Expensive– Once for all documents and queries (offline)
1 Query-independent – requires combination with query-dependent criteria
1 Hard to spam | <ul style="list-style-type: none">1 Computation:<ul style="list-style-type: none">– Expensive– Requires computation for each query
1 Query-dependent
1 Relatively easy to spam1 Quality depends on quality of start set1 Gives hubs as well as authorities |
|--|---|

Google™ *Open problems*

- 1 Compare performance of query-dependent and query-independent connectivity analysis
- 1 Exploit order of links on the page (see e.g. [CDGKRRRT'98],[DH'99])
- 1 Both Google and HITS compute principal eigenvector. What about non-principal eigenvector? ([K'98])
- 1 Derive other graphs from the hyperlink structure ...





Algorithmic issues related to search engines

1 Collecting documents

- Priority
- Load balancing
 - Internal
 - External
- Trap avoidance
- ...

1 Processing and representing the data

- Query-independent ranking
- Graph representation
- Index building
- Duplicate elimination
- Categorization
- ...

1 Processing queries

- Query-dependent ranking
- Duplicate elimination
- Query refinement
- Clustering
- ...

- 1 Graphs derived from the hyperlink structure of the Web:
 - **Node** =page
 - **Edge** (u,v) iff pages u and v are related in a specific way (directed or not)
- 1 Examples of edges:
 - iff u has hyperlink to v
 - iff there exists a page w pointing to both u and v
 - iff u is often retrieved within x seconds after v
 - ...

- 1 Ranking algorithms
 - PageRank
 - HITS
 - ...
- 1 Search-by-example [DH'99]
- 1 Categorization of Web pages
 - [CDI'98]
- 1 Visualization/Navigation
 - Mapuccino [MJSUZB'97]
 - WebCutter [MS'97]
 - ...
- 1 Structured Web query tools
 - WebSQL [AMM'97]
 - WebQuery [CK'97]
 - ...



Example: SRC Connectivity Server [BBHKV'98]

Directed edges = Hyperlinks

1 **Goal:** Support two basic operations for all URLs collected by AltaVista

– InEdges(URL u , int k)

- Return k URLs pointing to u

– OutEdges(URL u , int k)

- Return k URLs that u points to

1 **Difficulties:**

– Memory usage (~180 M nodes, 1B edges)

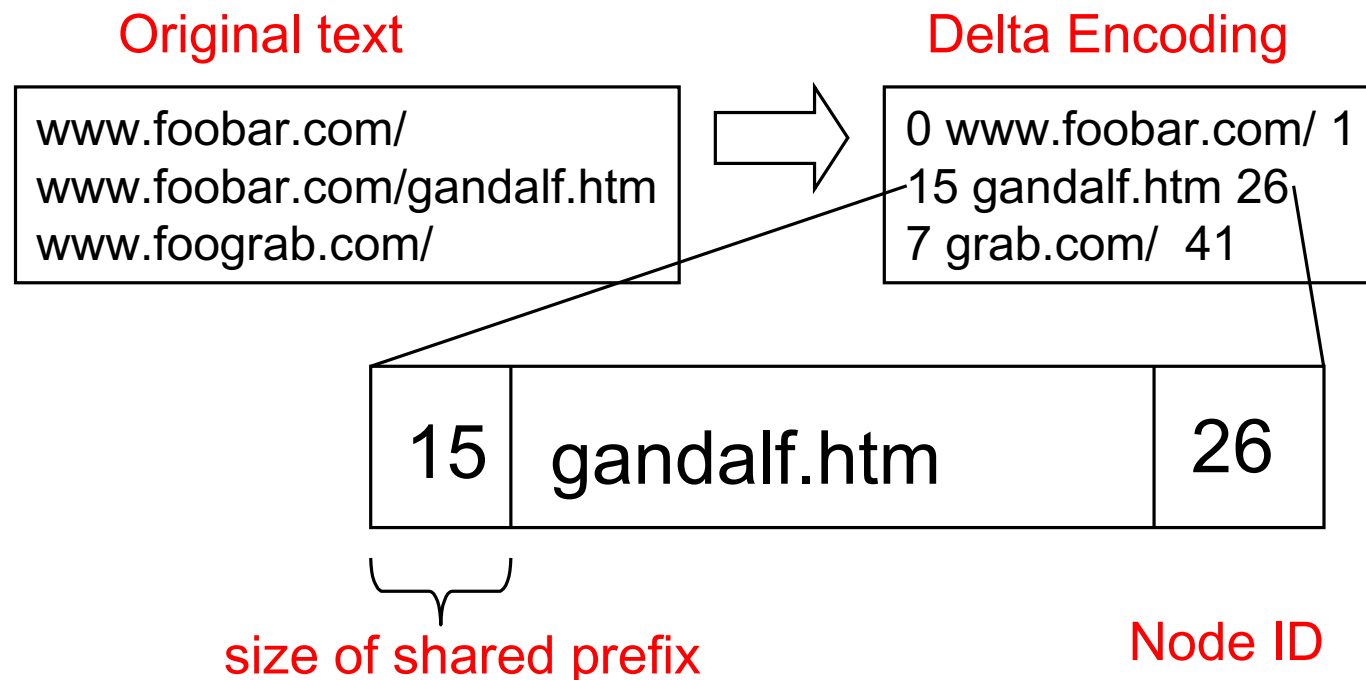
– Preprocessing time (days ...)

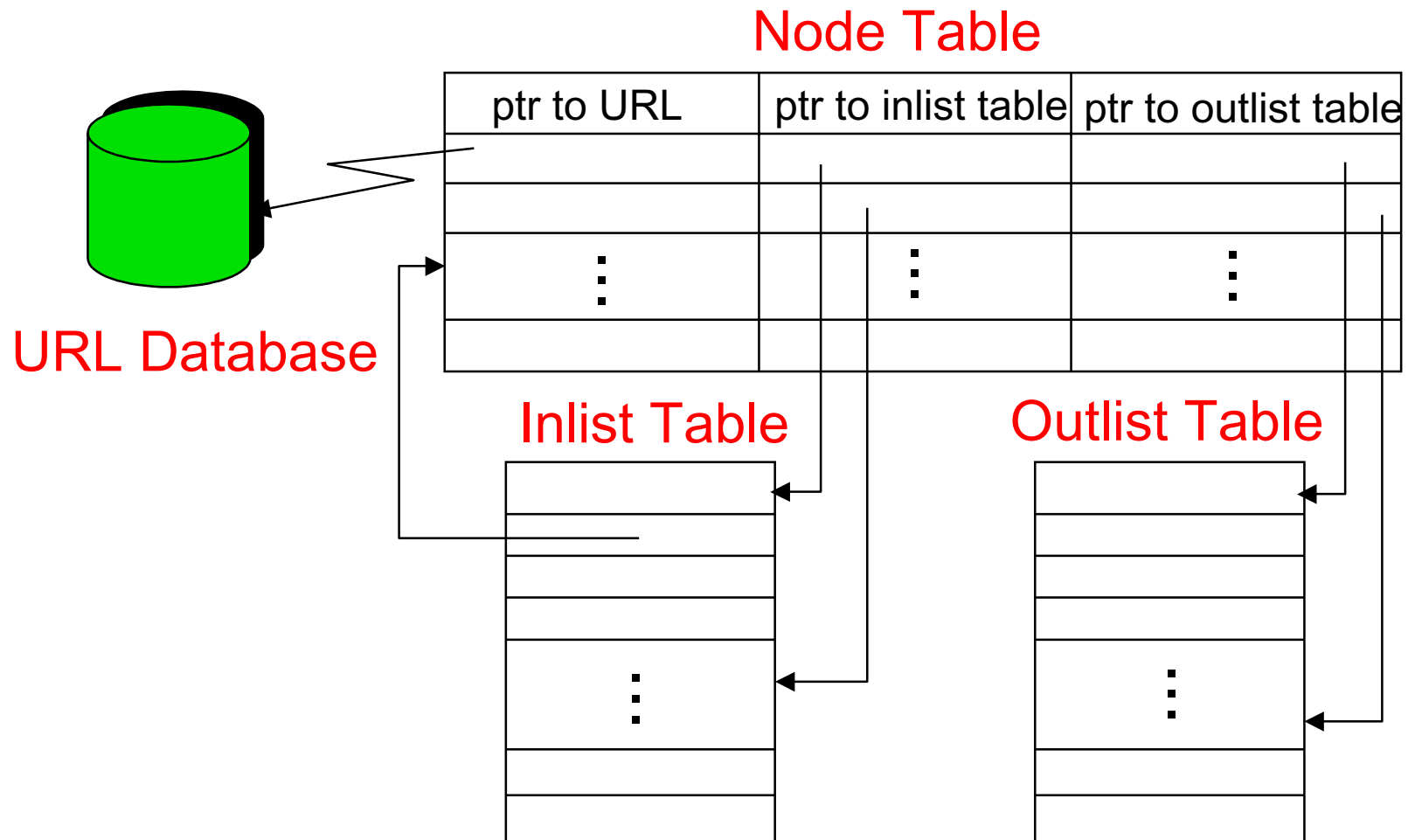
– Query time (~ 0.0001s/result URL)

Google *URL database*

Sorted list of URLs is 8.7 GB (\approx 48 bytes/URL)

Delta encoding reduces it to 3.8 GB (\approx 21 bytes/URL)





- 1 >1B nodes (12/99), mean indegree is ~ 8
- 1 *Zipfian Degree Distributions* [KRRT'99]:
 - $F_{in}(i)$ = fraction of pages with indegree i

$$F_{in}(i) \sim \frac{1}{i^{2.1}}$$

- $F_{out}(i)$ = fraction of pages with outdegree i

$$F_{out}(i) \sim \frac{1}{i^{2.38}}$$

Google™ *Open problems*



- 1 **Graph compression:** How much compression possible without significant run-time penalty?
 - Efficient algorithms to find frequently repeated small structures (e.g. wheels, $K_{2,2}$)
- 1 **External memory graph algorithms:** How to assign the graph representation to pages so as to reduce paging? (see [NGV'96, AAMVV'98])
- 1 **Stringology:** Less space for URL database? Faster algorithms for URL to node translation?
- 1 **Dynamic data structures:** How to make updates efficient at the same space cost?



Algorithmic issues related to search engines

1 Collecting documents

- Priority
- Load balancing
 - Internal
 - External
- Trap avoidance
- ...

1 Processing and representing the data

- Query-independent ranking
- Graph representation
- Index building
- Duplicate elimination
- Categorization
- ...

1 Processing queries

- Query-dependent ranking
- Duplicate elimination
- Query refinement
- Clustering
- ...

Google *Index building*

- 1 ***Inverted index data structure***: Consider all documents concatenated into one huge document
 - For each word keep an ordered array of all positions in document, potentially compressed

Word 1	1 st position	...	last position
⋮	⋮		⋮

- 1 Allows efficient implementation of AND, OR, and AND NOT operations



Algorithmic issues related to search engines

1 Collecting documents

- Priority
- Load balancing
 - Internal
 - External
- Trap avoidance
- ...

1 Processing and representing the data

- Query-independent ranking
- Graph representation
- Index building
- Duplicate elimination
- Categorization
- ...

1 Processing queries

- Query-dependent ranking
- Duplicate elimination
- Query refinement
- Clustering
- ...

- 1 Proliferation of almost equal documents on the Web:
 - Legitimate: Mirrors, local copies, updates, etc.
 - Malicious: Spammers, spider traps, dynamic URLs
 - Mistaken: Spider errors

- 1 Approximately 30% of the pages on the Web are (near) duplicates. [BGMZ'97,SG'98]



Uses of duplicate information

- 1 Smarter crawlers
- 1 Smarter web proxies
 - Better caching
 - Handling broken links
- 1 Smarter search engines
 - no duplicate answers
 - smarter connectivity analysis
 - less RAM and disk

Google™ *2 Types of duplicate filtering*

- 1 **Fine-grain:** Finding near-duplicate documents
- 1 **Coarse-grain:** Finding near-duplicate hosts (*mirrors*)



Fine-grain: Basic mechanism

- 1 Must filter both **duplicate** and **near-duplicate** documents
- 1 Computing pair-wise edit distance would take forever
- 1 Preferably to store only a **short** sketch for each document.



The basics of a solution

[B'97],[BGMZ'97]

1. Reduce the problem to a set intersection problem
2. Estimate intersections by sampling minima.

Google™ Shingling

- 1 Shingle = Fixed size sequence of w contiguous words

a rose is a rose is a rose

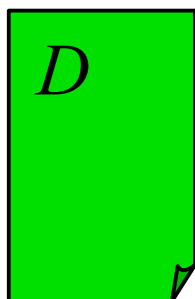
a rose is a

rose is a rose

is a rose is

a rose is a

rose is a rose

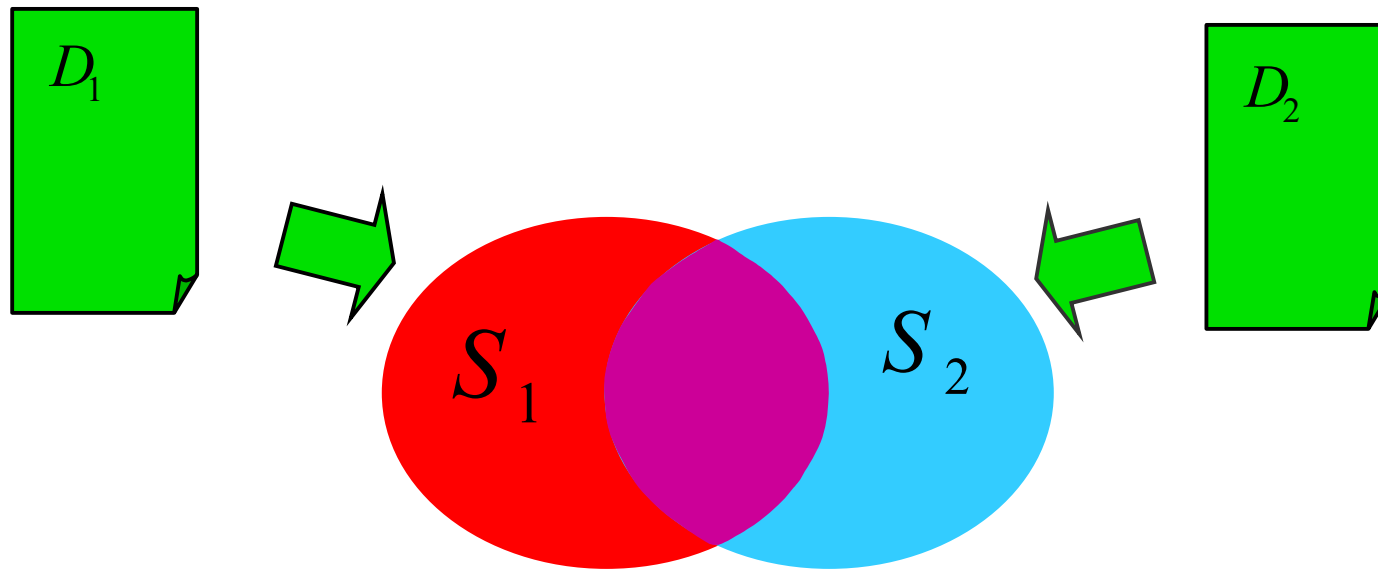


Shingling

Set of
shingles

Fingerprint

Set of
64 bit
fingerprints

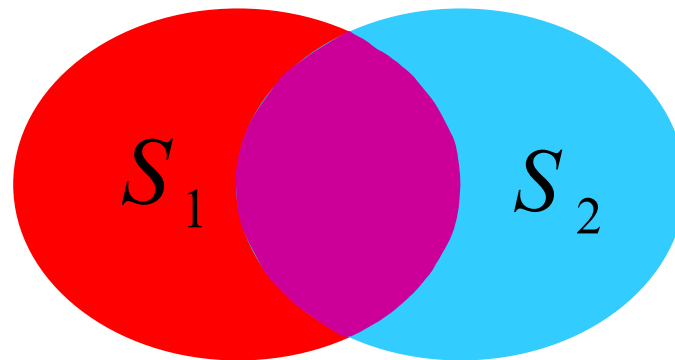


$$\text{resemblance} = \frac{|S_1 \cap S_2|}{|S_1 \cup S_2|}$$

Google *Sampling minima*

- 1 Apply a random permutation σ to the set $[0..2^{64}]$
- 1 Crucial fact

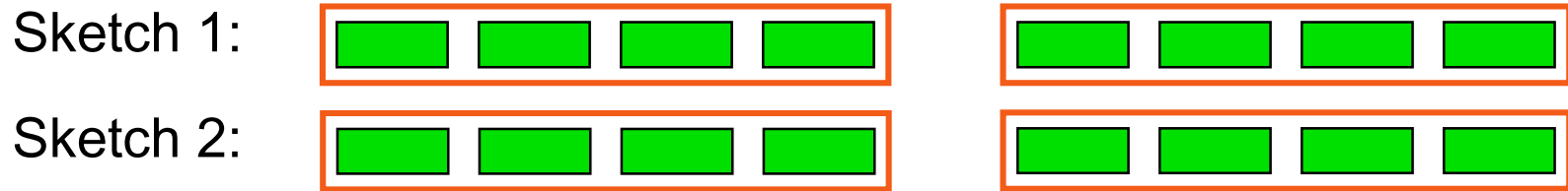
Let $\alpha = \min(\sigma(S_1))$ $\beta = \min(\sigma(S_2))$



$$\Pr(\alpha = \beta) = \frac{|S_1 \cap S_2|}{|S_1 \cup S_2|}$$

Google *Implementation*

- 1 Choose a set of t random permutations of U
- 1 For each document keep a sketch $S(D)$ consisting of t minima = **samples**
- 1 Estimate resemblance of A and B by counting common samples
- 1 The permutations should be from a **min-wise independent** family of permutations. See [BCFM'97] for the theory of **mw**i permutations.



- 1 Divide sketch into k groups of s samples ($t = k * s$)
- 1 Fingerprint each group \Rightarrow feature
- 1 Two documents are fungible if they have at least r common features.
- 1 Want
Fungibility \Leftrightarrow Resemblance above fixed threshold ρ

- 1 $\rho = 90\%$. In a 1000 word page with shingle length = 8 this corresponds to
 - Delete a paragraph of about 50-60 words.
 - Change 5-6 random words.
- 1 Sketch size $t = 84$, divide into $k = 6$ groups of $s = 14$ samples
- 1 8 bytes fingerprints \rightarrow we store only $6 \times 8 =$
48 bytes/document
- 1 Threshold $r = 2$

Probability that two documents are deemed fungible

Two documents with resemblance ρ

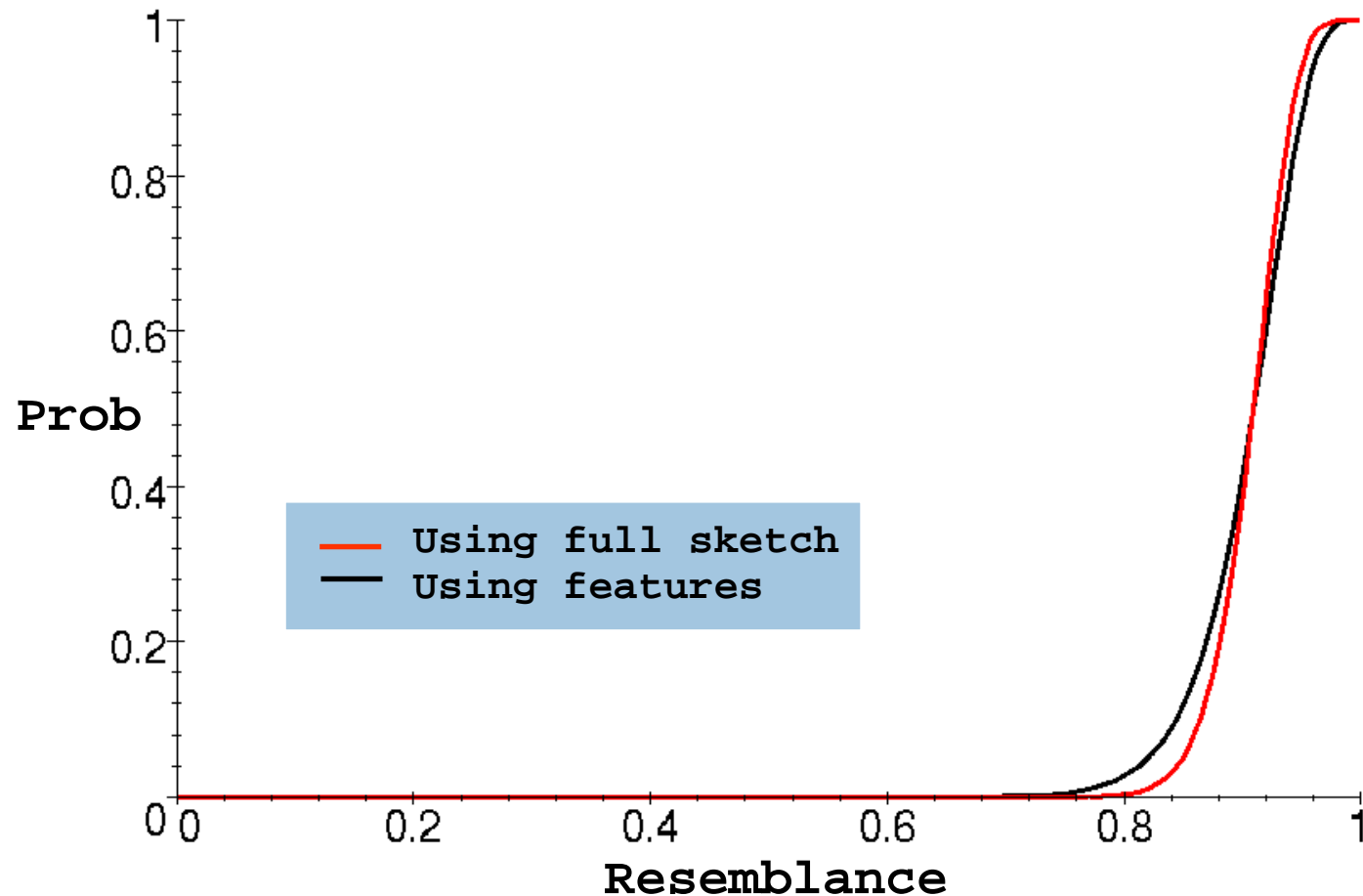
1 Using the full sketch

$$P = \sum_{i=r \cdot s}^{k \cdot s} \binom{k \cdot s}{i} \rho^i (1 - \rho)^{k \cdot s - i}$$

1 Using features

$$P = \sum_{i=r}^k \binom{k}{i} \rho^{s \cdot i} (1 - \rho^s)^{k - i}$$

Probability that two pages are deemed fungible



Fine-grain duplicate elimination: open problems and related work



- 1 Best way of grouping samples for a given threshold and/or for multiple thresholds?
- 1 Efficient ways to find in a data base pairs of records that share many attributes. Best approach?
- 1 Min-wise independent permutations -- lots of open questions.
- 1 Other applications possible (images, sounds, ...) -- need translation into set intersection problem.

- 1 Related work: M'94, BDG'95, SG'95, H'96, FSGMU'98

Google™ *2 Types of duplicate filtering*

Fine-grain: Finding near-duplicate documents

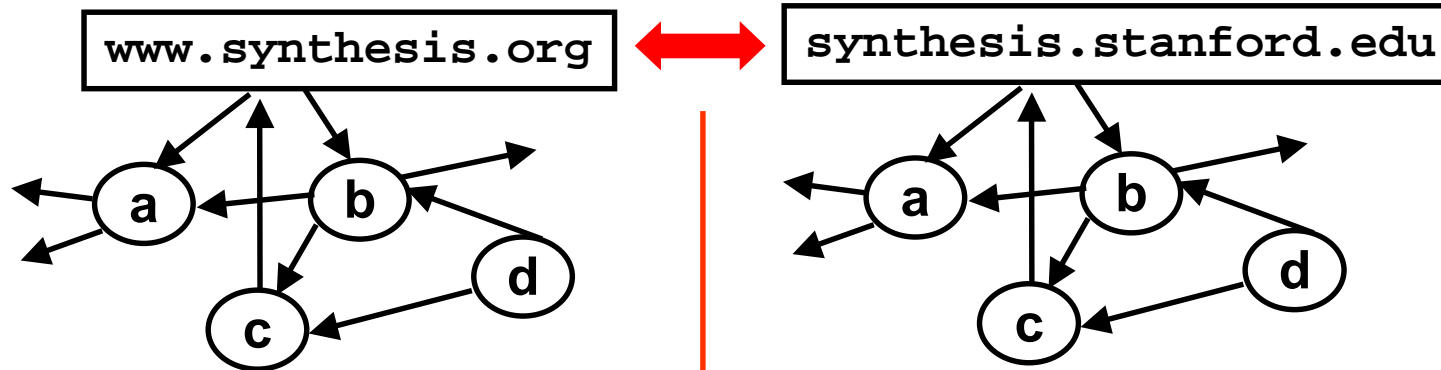
1 **Coarse-grain:** Finding near-duplicate hosts (*mirrors*)

1 Input:

- **Subset** of URLs on various hosts, collected e.g. by search engine crawl or web proxy
- No **content** of pages pointed to by URLs except each page is labeled with its out-links

1 Goal: Find pairs of hosts that mirror content

Google Example



www.synthesis.org/Docs/ProjAbs/synsys/synanalysis.html
www.synthesis.org/Docs/ProjAbs/synsys/visual-semi-quant.html
www.synthesis.org/Docs/annual.report96.final.html
www.synthesis.org/Docs/cicee-berlin-paper.html
www.synthesis.org/Docs/myr5
www.synthesis.org/Docs/myr5/cicee/bridge-gap.html
www.synthesis.org/Docs/myr5/cs/cs-meta.html
www.synthesis.org/Docs/myr5/mech/mech-intro-mechatron.html
www.synthesis.org/Docs/myr5/mech/mech-take-home.html
www.synthesis.org/Docs/myr5/synsys/experiential-learning.html
www.synthesis.org/Docs/myr5/synsys/mm-mech-dissec.html
www.synthesis.org/Docs/yr5ar
www.synthesis.org/Docs/yr5ar/assess
www.synthesis.org/Docs/yr5ar/cicee
www.synthesis.org/Docs/yr5ar/cicee/bridge-gap.html
www.synthesis.org/Docs/yr5ar/cicee/comp-integ-analysis.html

synthesis.stanford.edu/Docs/ProjAbs/deliv/high-tech-classroom.html
synthesis.stanford.edu/Docs/ProjAbs/mech/mech-enhanced-circ-anal.html
synthesis.stanford.edu/Docs/ProjAbs/mech/mech-intro-mechatron.html
synthesis.stanford.edu/Docs/ProjAbs/mech/mech-mm-case-studies.html
synthesis.stanford.edu/Docs/ProjAbs/synsys/quant-dev-new-teach.html
synthesis.stanford.edu/Docs/annual.report96.final.html
synthesis.stanford.edu/Docs/annual.report96.final_fn.html
synthesis.stanford.edu/Docs/myr5/assessment
synthesis.stanford.edu/Docs/myr5/assessment/assessment-main.html
synthesis.stanford.edu/Docs/myr5/assessment/mm-forum-kiosk-A6-E25.html
synthesis.stanford.edu/Docs/myr5/assessment/neato-ucb.html
synthesis.stanford.edu/Docs/myr5/assessment/not-available.html
synthesis.stanford.edu/Docs/myr5/cicee
synthesis.stanford.edu/Docs/myr5/cicee/bridge-gap.html
synthesis.stanford.edu/Docs/myr5/cicee/cicee-main.html
synthesis.stanford.edu/Docs/myr5/cicee/comp-integ-analysis.html



Coarse-grain: Basic mechanism

- 1 Must filter both **duplicate** and **near-duplicate** mirrors
- 1 Pair-wise testing would take forever
- 1 Both high precision (not outputting wrong mirrors) and high recall (finding almost all mirrors) are important



A definition of mirroring

Host1 and Host2 are mirrors iff

For all paths p such that

`http://Host1/p`

is a web page,

`http://Host2/p`

exists with duplicate (or *near-duplicate*) content,
and vice versa.



The basics of a solution

[BBDH'99]

1. Pre-filter to create a small set of pairs of potential mirrors (*pre-filtering step*)
2. Test each pair of potential mirrors (*testing step*)
3. Use different pre-filtering algorithms to improve recall

Google™ *Testing step*

- 1 Test root pages + x URLs from each host sample
- 1 **If** one test returns “not near-duplicate”
then hosts are *not mirrors*
- 1 **If** root pages and $> c\% x$ URLs from each host sample are near-identical
then hosts are *mirrors*,
else they are *not mirrors*

Google™ *Pre-filtering step*

- 1 Goal: Output **quickly** list of pairs of potential mirrors containing
 - many true mirror pairs (high recall)
 - not many non-mirror pairs (high precision)
- 1 Note: 2-sided error is allowed
 - Type-1: true mirror pairs might be missing in output
 - Type-2: non-mirror pair might be output
- 1 Testing of host pairs will eliminate type-2 errors, but not type-1 errors



Different pre-filtering techniques

- 1 IP-based
- 1 URL-string based
- 1 URL-string and hyperlink based
- 1 Hostname and hyperlink based

Problem with IP addresses

The image shows a Netscape browser window with two tabs. The top tab is titled "Factory - Netscape" and has the address bar set to `http://eliza-iii.ibex.co.nz/`. The bottom tab is titled "Factory - Netscape" and has the address bar set to `http://pixel.ibex.co.nz/`. A yellow box labeled `203.29.170.23` has arrows pointing to the domain names in both browser tabs. A terminal window titled "xterm" is overlaid on the browser, showing the following output:

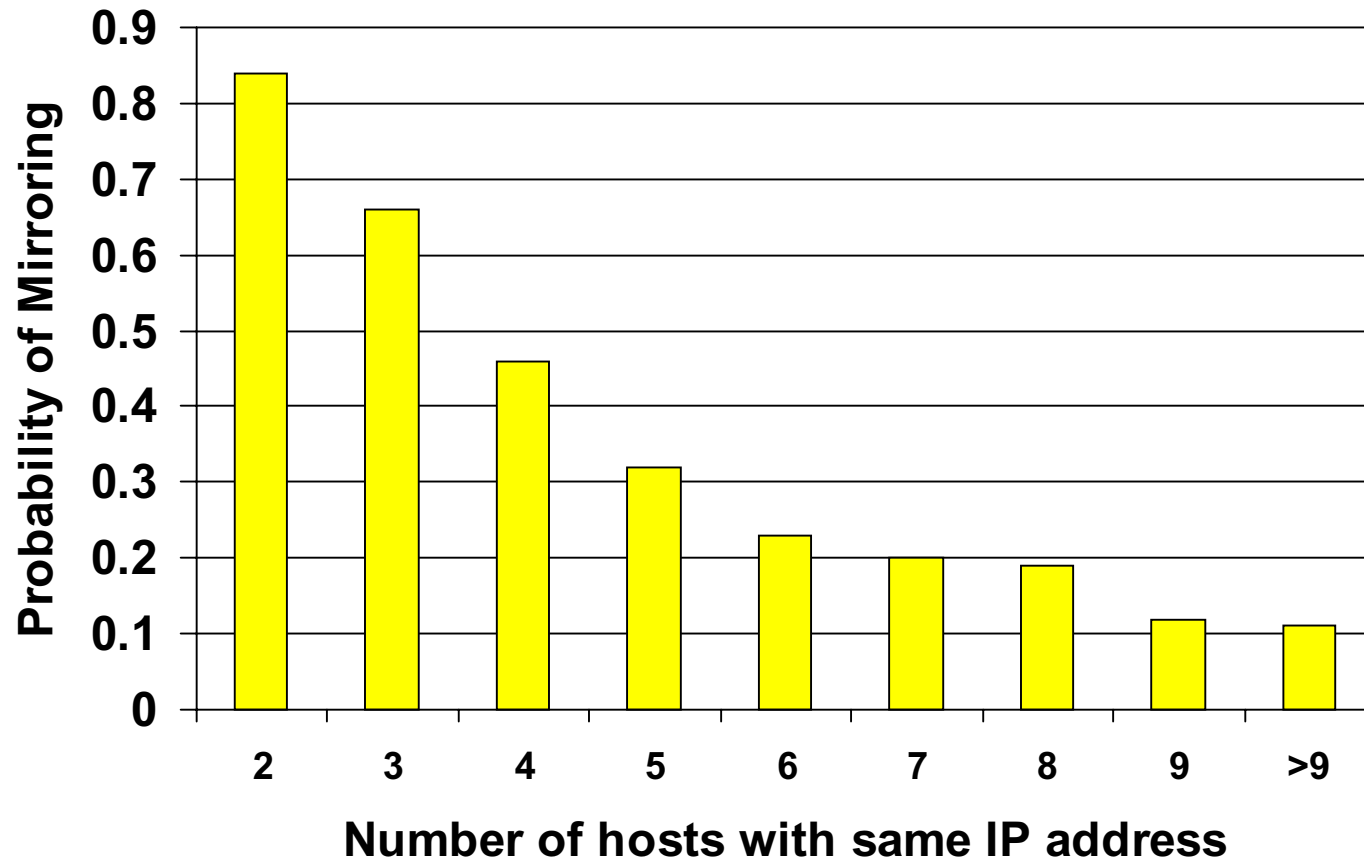
```
ash 1> nslookup eliza-iii.ibex.co.nz
Server: snakeoil.pa.dec.com
Address: 16.4.0.250
Name: eliza-iii.ibex.co.nz
Address: 203.29.170.23

ash 2> nslookup pixel.ibex.co.nz
Server: snakeoil.pa.dec.com
Address: 16.4.0.250
Name: pixel.ibex.co.nz
Address: 203.29.170.23

ash 3>
```

The browser window also displays a "Pixel Factory" logo and navigation links: [Main Page](#), [Comments](#), [Links](#), and [Mail Chris](#).

Number of host with same IP address vs mirror probability





IP based pre-filtering algorithms

- 1 *IP4*: Cluster hosts based on IP address
 - Enumerate pairs from clusters in increasing cluster size (max 200 pairs)

- 1 *IP3*: Cluster hosts based on first 3 octets of their IP address
 - Enumerate pairs from clusters in increasing cluster size (max 5 pairs)



URL string based pre-filtering algorithms

Information extracted from URL strings:

- 1 Similar hostnames: might belong to same organization
- 1 Similar paths: might have replicated directories
- 1 extract “features” for host from URL strings and test similarity

Similarity Testing Approach:

- 1 *Feature vector* for each host similar to term vector for document:
 - Host corresponds to document
 - Feature corresponds to term
- 1 Similarity of hosts = Cosine of angle of feature vectors



URL string based algorithms (cont.)

- *paths*: Features are paths: e.g.,
/staff/homepages/~dilbert/foo
- *prefixes*: Features are prefixes: e.g.,
/staff
/staff/homepages
/staff/homepages/~dilbert
/staff/homepages/~dilbert/foo
- Other variants: *hosts* and *shingles*

Google™ *Paths + connectivity (conn)*

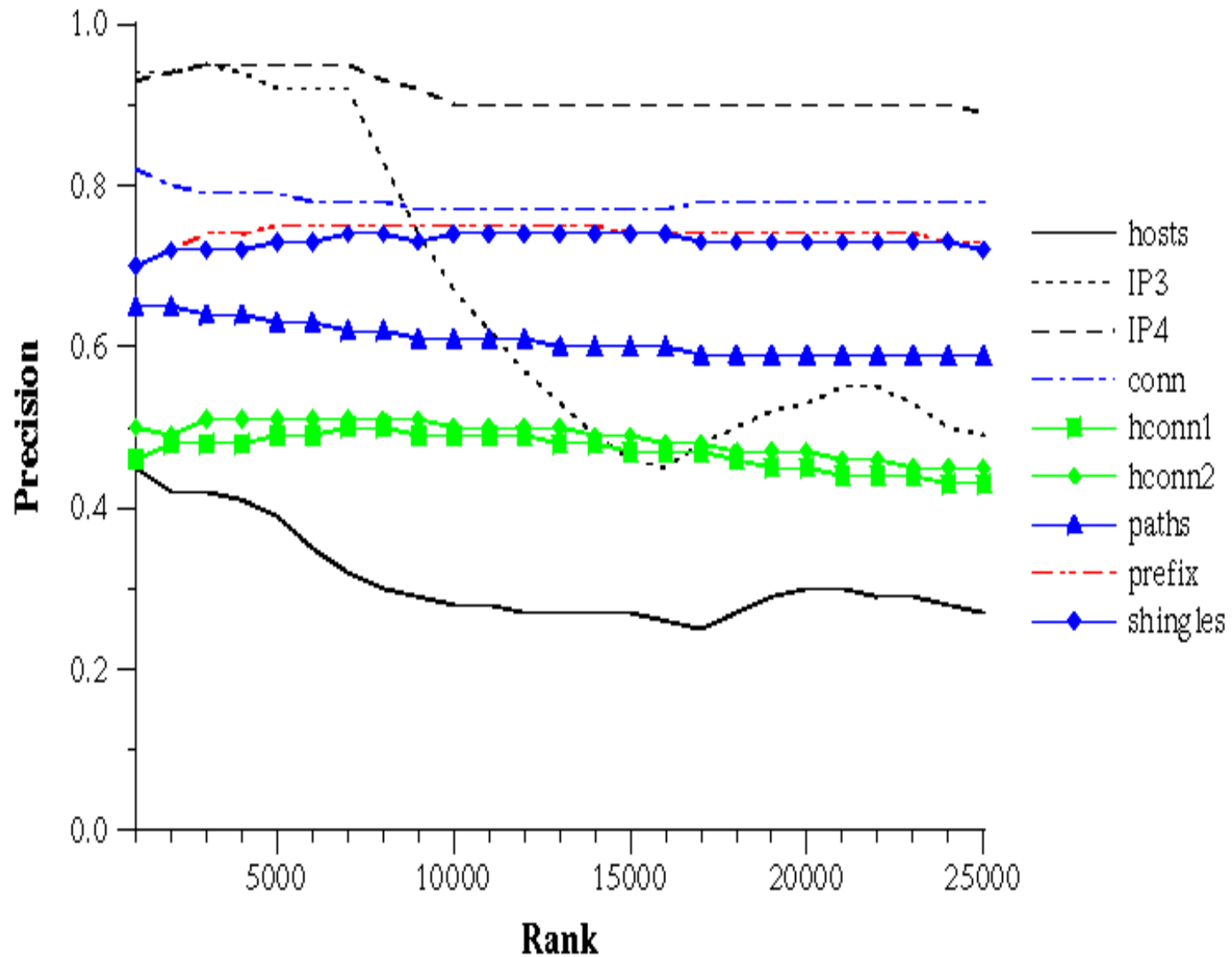
- 1 Take output from *paths* and filter thus:
 - Consider 10 common paths in sample with highest outdegree
 - Paths are **equivalent** if 90% of their combined out-edges are common to both
 - Keep host-pair if 75% of the paths are equivalent

- 1 **Idea:** Mirrors point to similar set of other hosts
- 2 Feature vector approach to test similarity:
 - features are hosts that are pointed to
 - 2 different ways of feature weighting:
 - *hconn1*
 - *hconn2*

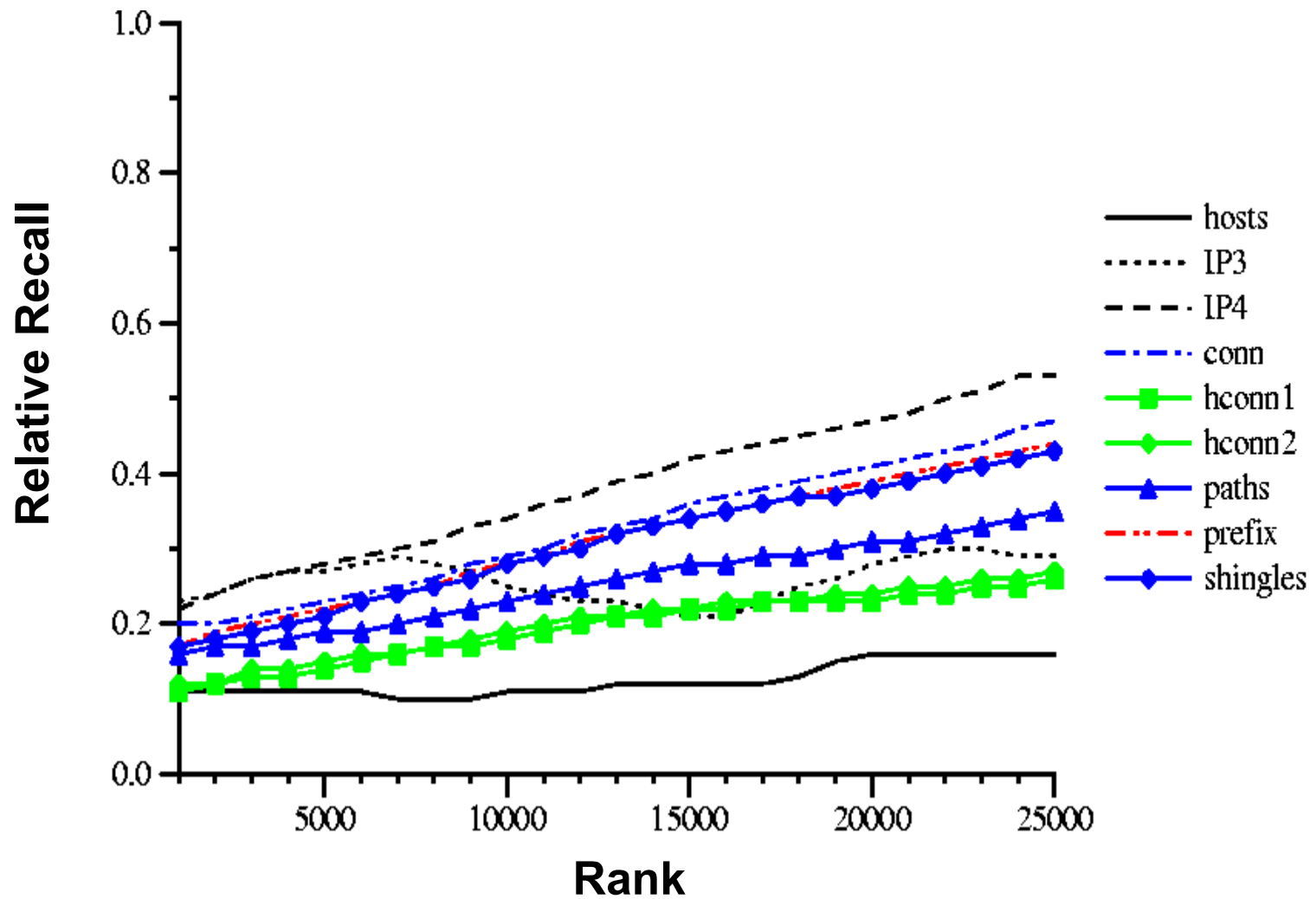
- 1 Input: 140 million URLs on 233,035 hosts + out-edges
 - Original 179 million URLs reduced by considering only hosts with at least 100 URLs in set
- 1 For each of the above pre-filtering algorithms:
 - Compute list of 25,000 (100,000) ranked pairs of potential mirrors
 - Test each pair of potential mirrors (testing step) and output list of mirrors

Determine precision and *relative recall*

Precision up to rank 25,000



Relative recall up to rank 25,000





Relative recall at 25,000 for combined output

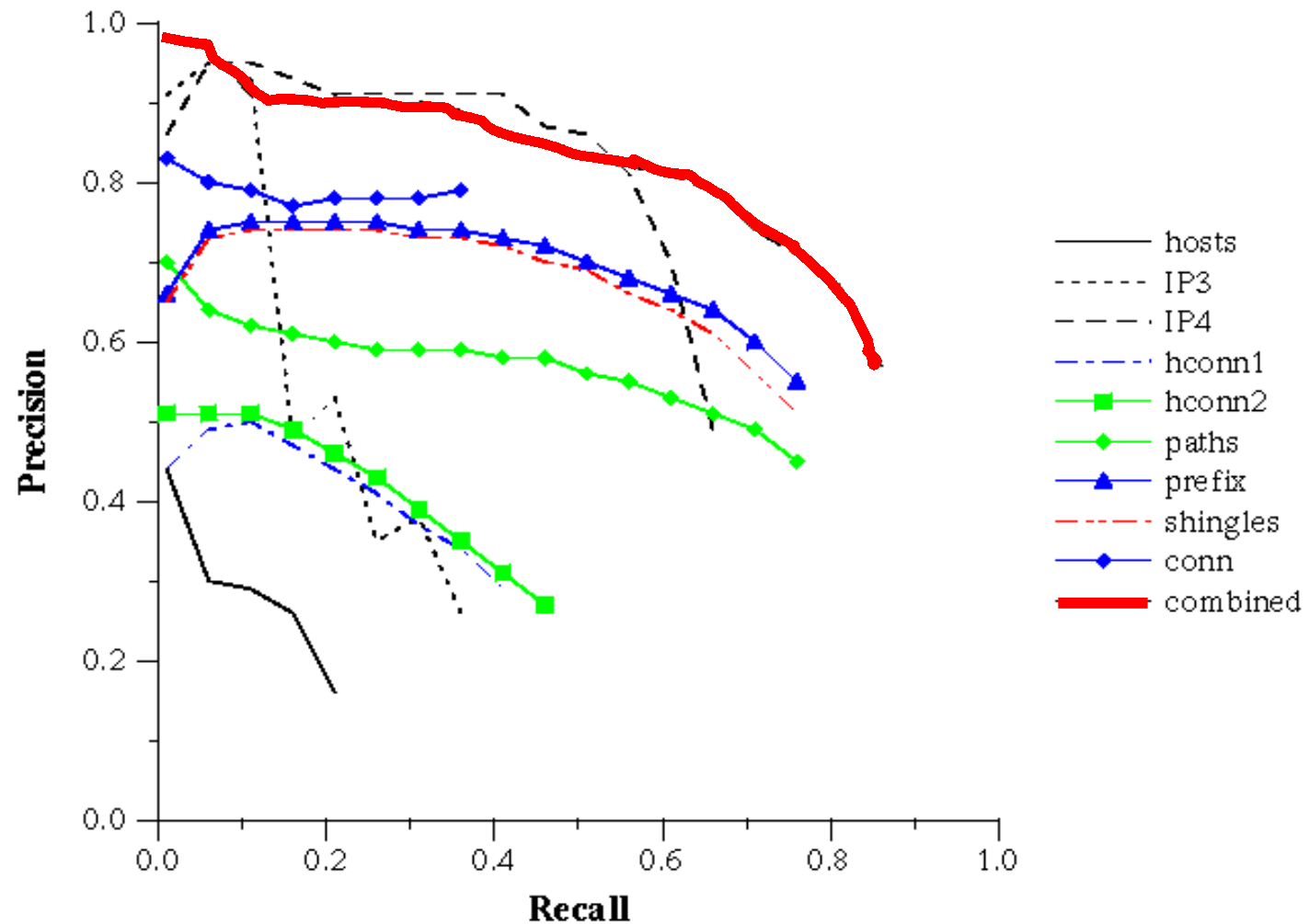
	<i>hosts</i>	<i>IP3</i>	<i>IP4</i>	<i>conn</i>	<i>hconn1</i>	<i>hconn2</i>	<i>paths</i>	<i>prefix</i>	<i>shingles</i>
<i>hosts</i>	17%								
<i>IP3</i>	39%	30%							
<i>IP4</i>	61%	58%	54%						
<i>conn</i>	58%	66%	80%	47%					
<i>hconn1</i>	40%	51%	69%	59%	26%				
<i>hconn2</i>	41%	52%	70%	60%	29%	27%			
<i>paths</i>	48%	59%	78%	55%	51%	52%	36%		
<i>prefix</i>	54%	61%	75%	65%	58%	58%	57%	44%	
<i>shingles</i>	53%	61%	75%	64%	57%	58%	57%	48%	44%



Combined approach (combined)

- 1 Combines top 100,000 results from *hosts*, *IP4*, *paths*, *prefix*, and *hconn1*.
- 1 Sort host pairs by:
 - Number of algorithms that return the host pair
 - Use best rank for any algorithm to break ties between host pairs
- 1 At rank 100,000: relative recall of 86%, precision of 57%

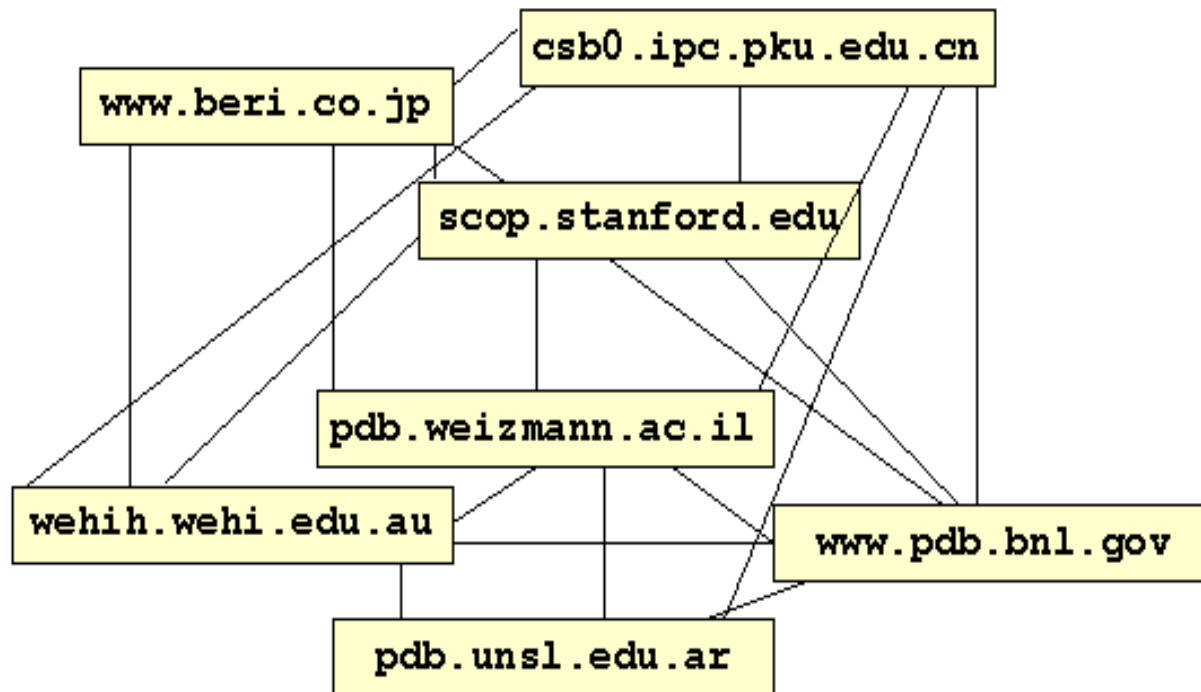
Precision vs relative recall



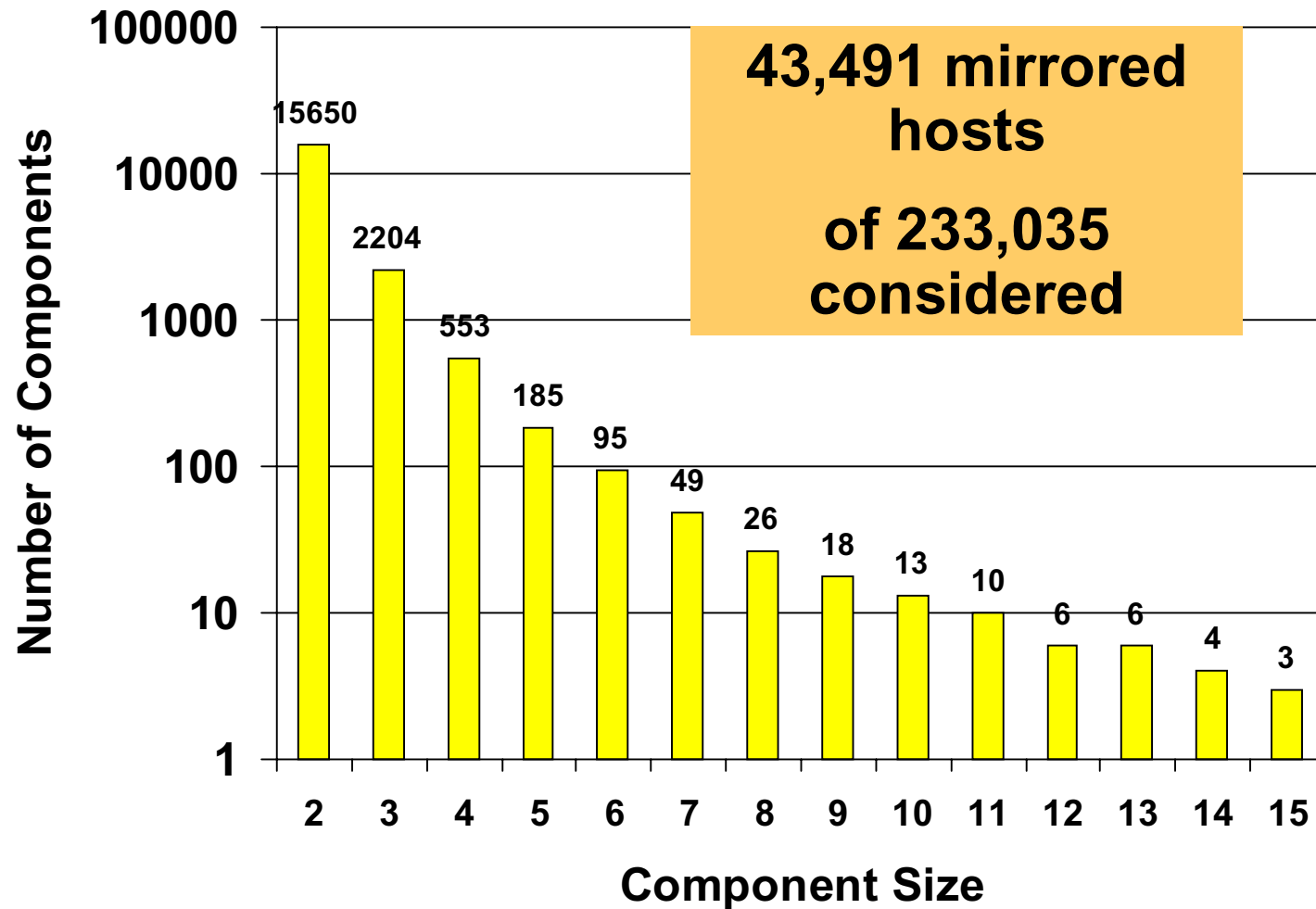
Google™ *Web host graph*

- 1 A **node** for each host h
- 1 An **undirected edge** (h, h') if h and h' are output as mirrors
- ⊞ Each **(connected) component** gives a set of mirrors

Protein Data Bank



Component size distribution





Coarse-grain duplicate filtering: Summary and open problems

- 1 Mirroring is common (43,491 mirrored hosts out of 233,035 considered hosts)
 - Load balancing, franchises/branding, virtual hosting, spam
- 1 Mirror detection based on non-content attributes is feasible.
- 1 [CSG'00] use page content similarity based approach.
Open Problem: Compare and combine content and non-content techniques.
- 1 Open Problem: Assume you can choose which URLs to visit at a host. Determine best technique.



Algorithmic issues related to search engines

1 Collecting documents

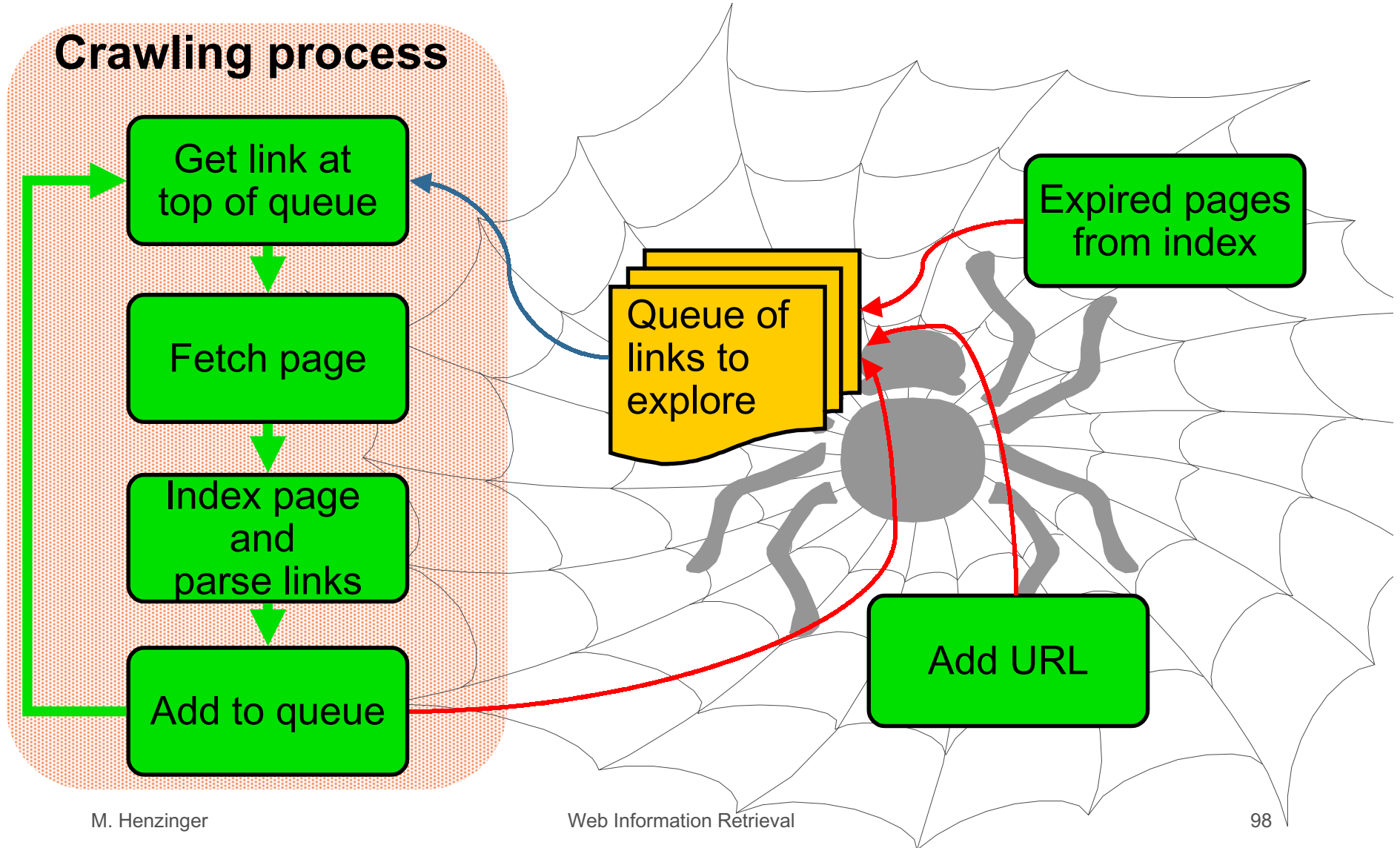
- Priority
- Load balancing
 - Internal
 - External
- Trap avoidance
- ...

1 Processing and representing the data

- Query-independent ranking
- Graph representation
- Index building
- Duplicate elimination
- Categorization
- ...

1 Processing queries

- Query-dependent ranking
- Duplicate elimination
- Query refinement
- Clustering
- ...



Google™ *Queuing discipline*

- 1 Standard graph exploration:
 - Random
 - BFS
 - DFS (+ depth limits)
- 1 **Goal:** get “best” pages for a given index size
 - Priority based on query-independent ranking:
 - highest indegree [M’95]
 - highest potential PageRank [CGP’98]
- 1 **Goal:** keep index fresh
 - Priority based on rate of change [CLW’97]



Google™ *Load balancing*

1 **Internal** -- can not handle too much retrieved data simultaneously, but



- Response time is unpredictable
- Size of answers is unpredictable
- There are additional system constraints (# threads, # open connections, etc.)

1 **External**



- Should not overload any server or connection
- A well-connected crawler can saturate the entire outside bandwidth of some small countries
- Any queuing discipline must be acceptable to the community

General-purpose search engines

- 1 Hierarchical directories
- 1 Specialized search engines
(dealing with heterogeneous data sources)
- 1 Search-by-example
- 1 Collaborative filtering
- 1 Meta-information

Building of hierarchical directories:

1 **Manual:** Yahoo!, LookSmart, Open Directory

1 **Automatic:**

- Populating of hierarchy [CDRRGK'98]: For each node in the hierarchy formulate fine-tuned query and run modified HITS algorithm
- Categorization: For each document find “best” placement in the hierarchy. Techniques are connectivity and/or text based [CDI'98, ...]

General-purpose search engines

Hierarchical directories

- 1 Specialized search engines
(dealing with heterogeneous data sources)
 - Shopping robots
 - Home page finder [SLE'97]
 - Applet finders
 - ...
- 1 Search-by-example
- 1 Collaborative filtering
- 1 Meta-information



Dealing with heterogeneous sources

1 Modern life problem:

Given information sources with various capabilities, query all of them and combine the output.

1 Examples

- Inter-business e-commerce e.g. `www.industry.net`
- Meta search engines
- Shopping robots

1 Issues

- Determining relevant sources -- the “**identification**” problem
- Merging the results -- the “**fusion**” problem



Example: a shopping robot

- 1 **Input:** A product description in some form
- Find:** Merchants for that product on the Web
- 1 **Jango** [DEW'97]
 - **preprocessing:** Store vendor URLs in database; learn for each vendor:
 - the URL of the search form
 - how to fill in the search form and
 - how the answer is returned
 - **request processing:** fill out form at every vendor and test whether the result is a success
 - range of products is predetermined

Excite Product Finder

powered by **Jango**

Help!

Need an example?
Try this:
Enter "Chardonnay"
for Variety and
"1991" for Year.
Then click "Find
Prices" or "Find
Reviews."

Still not sure what to
do? [Click here](#) for
further instructions.

**More Excite
Links**

[More Food & Drink](#)

Find Product Prices & Reviews

Know what you're shopping for? Find product information fast by entering at least one detail in the form below and clicking "Find Prices" or "Find Reviews." For a different selection of products in this category, click one of the links to the right.

Wine

Variety:

Winery:

Year:

Find Prices
Find Reviews

**Gourmet &
Groceries
Categories**

[Coffee](#)

[Tea](#)

> [Wine](#)

[Spirits](#)

[Liqueurs](#)

Excite Product Finder

powered by **Jango**

Wine - Products


[Get Reviews](#) [New Wine Search](#)

Your Search: Variety = "viognier"

Instructions: **Click a column title to sort results** by the information in that column. For more details on a particular wine, click on a link in the Name column.

Search Results: 15 items have been located. [Click here](#) for a search summary.

<u>Winery</u>	<u>Variety</u>	<u>Name</u>	<u>Year</u>	<u>Quantity</u>	<u>Store</u>	<u>Price</u>	
	Viognier	Alban Vineyards Estate Viognier	97	1 bottle	K&L Wine Merchants	\$22.95	Buy!
Arrowood Vineyards & Winery	Viognier	Arrowood Viognier Saralee's Vineyard	96	1 bottle	California Wine	\$28.00	Buy!
Calera	Viognier	Calera Mt. Harlan	95	1 bottle	Taylor & Norton	\$26.99	Buy!
Calera	Viognier	Calera Mt. Harlan	97	1 bottle	Taylor & Norton	\$27.49	Buy!
	Viognier	Chance Creek Viognier	97	1 bottle	K&L Wine Merchants	\$13.99	Buy!
	Viognier	Gregory Graham	97	1 bottle	Taylor & Norton	\$19.99	Buy!



"One of the nations top 10 wine retailers"
klwines.com -Publishers of
The Wine Spectator

Please enter the keyword by which you would like to search:

Search Fields: Item Name Description Both Fields

Search

Hint for searching... The search engine only returns EXACT matches. Use terms that are not overly specific. For example, instead of searching for "Beringer Private Reserve Cabernet Sauvignon" simply search for "Beringer"



Wines	Shopping Basket 🛒	Wine-of-the-Month	Contact Us
· Domestic Wine · Imported Wine · Old and Rare Wine	Port Malts Wine Info	Your Account Ordering Info	
	Wine Accessories	Search K&L 🔍	FREE Newsletter

[Copyright and disclaimer](#) © K&L Wine Merchants

- **1983 Pichon Lalande** **Qnty. Avail.: 3** **Bottle: \$129.00**
94 points from Parker... 'Consistently one of the great wines of the 1983 vintage, as well as one of my personal favorites, this beautiful wine has been gorgeous to drink since bottling. It displays no signs of evolution, although it remains undeniably rich, seductive, and compelling. Deep dark ruby-colored, with a huge nose of Asian spices, blackcurrants, plums, and flowers, this super-concentrated, velvety-textured wine reveals gobs of rich, creamy fruit. It can be drunk now or cellared for 15-20 years. It is Pauillac at its most **DECADENT** and seductive!'
- **1961 Palmer, Margaux** **Qnty. Avail.: 5** **Bottle: \$699.00**
99 points from Robert Parker... 'The 1961 Palmer has long been considered to be a legend from this vintage, and its reputation is well-deserved. The wine is at its apogee, with an extraordinary, sweet, complex nose with aromas of flowers, cassis, toast, and minerals. It is intensely concentrated, offering a cascade of lavishly ripe, full-bodied, opulent fruit, soft tannins, and a voluptuous finish. This is a **DECADENT** Palmer, unparalleled since in quality with the exception of 1983 and 1989.'
- **1996 Charmes-Chambertin, Jean Raphet** **Bottle: \$79.95**
Top Pick! The obvious standout in a super lineup of Raphet wines. Very tasty and very limited. 95 points from the Wine Advocate (Pierre Rovanni)... 'Extraordinary. Medium-to-dark rub-colored, its awesome aromatics reveal fruit cake, cinnamon, spicy chutney, and assorted red fruits. This is a sexy, intense, **DECADENT**, and full-bodied gem crammed with loads of sweet cherries, perfume, flowers and spices... the mouthwatering flavors continue to coat the palate for what seems like minutes.'
- **1993 Haut-Marbuzet, St-Estephe** **Bottle: \$19.95**
Ripe cherry-berry scents in the nose are followed by a delicious, lush, elegant wine that gently flows onto the palate. Hints of cedar and tobacco. Robert Parker... 'Haut-Marbuzet is one of the oldest estates in St.-Estephe, but its fame can be traced only to 1952, when it was purchased by the father of the current proprietor, Henri Duboscq. No one argues with the success proprietor Duboscq has enjoyed. His wine is a Bordeaux that behaves more like a **DECADENT** Burgundy or Rhone.'

Google™ *Open problems*



- 1 Going beyond the lowest common capability
- 1 Learning problem: automatic understanding of new interfaces
- 1 Efficient ways to determine which DB is relevant to a particular query
- 1 “Partial knowledge indexing”: indexer has limited access to the full DB

General-purpose search engines

Hierarchical directories

Specialized search engines

(dealing with heterogeneous data sources)

- 1 Search-by-example
- 1 Collaborative filtering
- 1 Meta-information

Google™ *Search-by-example*

- 1 **Given:** set S of URLs
- 1 **Find:** URLs of similar pages
- 1 Various Approaches:
 - **Connectivity-based**
 - **Usage based:** related pages are pages visited frequently after S
 - **Query refinement**
- 1 Related Work: G'72, G'79, K'98, PP'97, S'97, CDI'98



Output from Google: *related:www.ebay.com*

[Fannie Mae's HomePath.com - Your On-Line Path to a Home of Your Own](#)

Fannie Mae's consumer Web site provides comprehensive information on buying and refinancing a home. Homebuyers can find -- Welcome to HomePath &re; -- a site that will help yo...

[www.homepath.com/](#) [Cached \(7k\)](#) **New!** Try out [GoogleScout](#)

[Today's Mortgage Information from HSH Associates, Financial Publishers](#)

HSH Associates, the world's leading publisher of mortgage and consumer loan information, surveys current loan rates from 2,500 to 3,000 lenders throughout the US. We offer dai...

[www.hsh.com/](#) [Cached \(15k\)](#) **New!** Try out [GoogleScout](#)

[Countrywide Home Loans](#)

The nation's largest independent mortgage lender. Pre-qualify for your potential maximum loan amount based on current interest rates and loan products with the Home Loan ...

[www.countrywide.com/](#) [Cached \(7k\)](#) **New!** Try out [GoogleScout](#)

[Keystroke Loans](#)

Home loans at the best interest rates from the Web's leading mortgage loan broker - Keystroke.com -- July 8, 1999 Check out our Mortgage Rates for purchase and refinance quote...

[www.keystrokenet.com/](#) [Cached \(8k\)](#) **New!** Try out [GoogleScout](#)

[www.igualify.com/](#)

New! Try out [GoogleScout](#)



Output from Alexa: www.ebay.com



You are here: [Home](#) > What's Related

What's Related

by  Alexa

...to <http://www.eloan.com/>

1. [Online Mortgage](#)
2. [Countrywide Home Loans](#)
3. [American Finance On Line](#)
4. [Keystroke Loans](#)
5. [Capital Mortgage Services, Inc.](#)
6. [Business Week](#)
7. [Chase Manhattan Mortgage Corporation](#)
8. [Home Loans](#)
9. [HomeByNet Home Page](#)
10. [1003 LOAN APPLICATION - APPLICATION FORMS - Mortgage broker, loan, interest](#)
11. [Learn About Smart Browsing...](#)

[DH'99]

1 Algorithm Companion

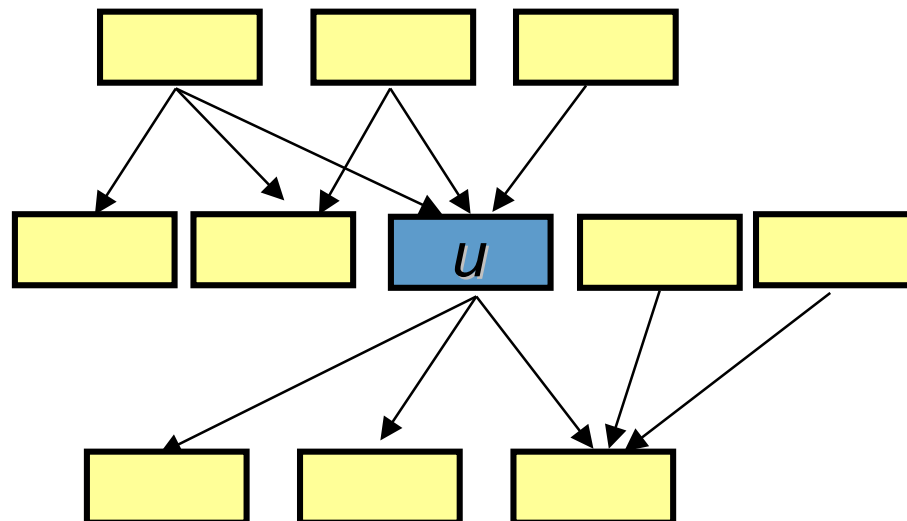
1 Algorithm Co-citation

Google *Algorithm Companion*

- 1 Build **modified neighborhood graph N** .
- 1 Run modified HITS algorithm on N .

Major Question: How to form neighborhood graph
s.t. top returned pages are useful related pages

- 1 **Node set:** From URL u go 'back', 'forward', 'back-forward', and 'forward-back'
- 1 **Edge set:** Directed edge if there is a hyperlink between 2 nodes
- 1 **Apply refinements** to N

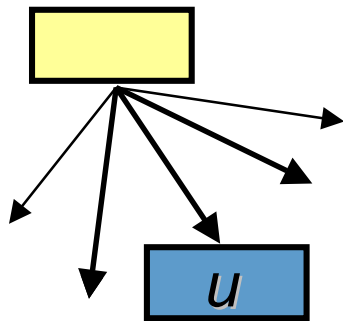


1 **Motivation:** Some nodes have high out-degree => graph would become too large

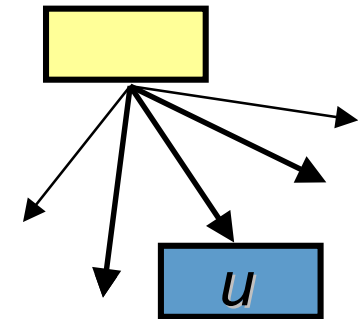
1 **Limit out-degree when going “forward”**

– Going forward from u : choose first 50 out-links on page

– Going forward from other nodes: choose 8 out-links **surrounding the in-link traversed to reach u**



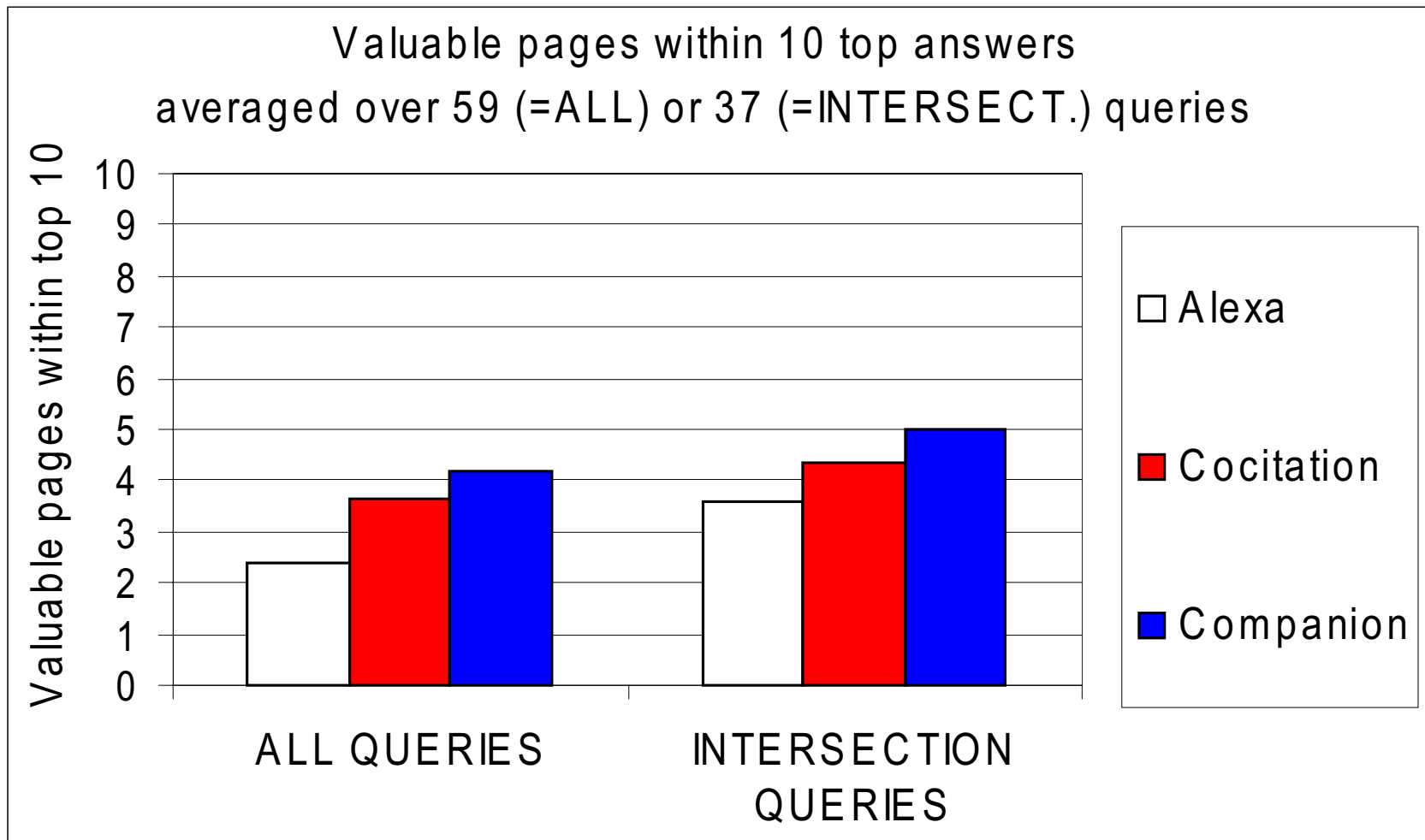
- 1 Determine 2000 arbitrary back-nodes of u .
- 1 Add to set S of siblings of u :
 - For each back node 8 forward-nodes surrounding the link to u
- 1 If there is enough co-citation with u then
 - return nodes in S in decreasing order of co-citations
- else
 - restart algorithm with one path element removed (http://.../X/Y/ -> http://.../X/)



1 Uses:

- Document Content
- Usage Data
- Connectivity

1 Removes path elements if no answer for u is found



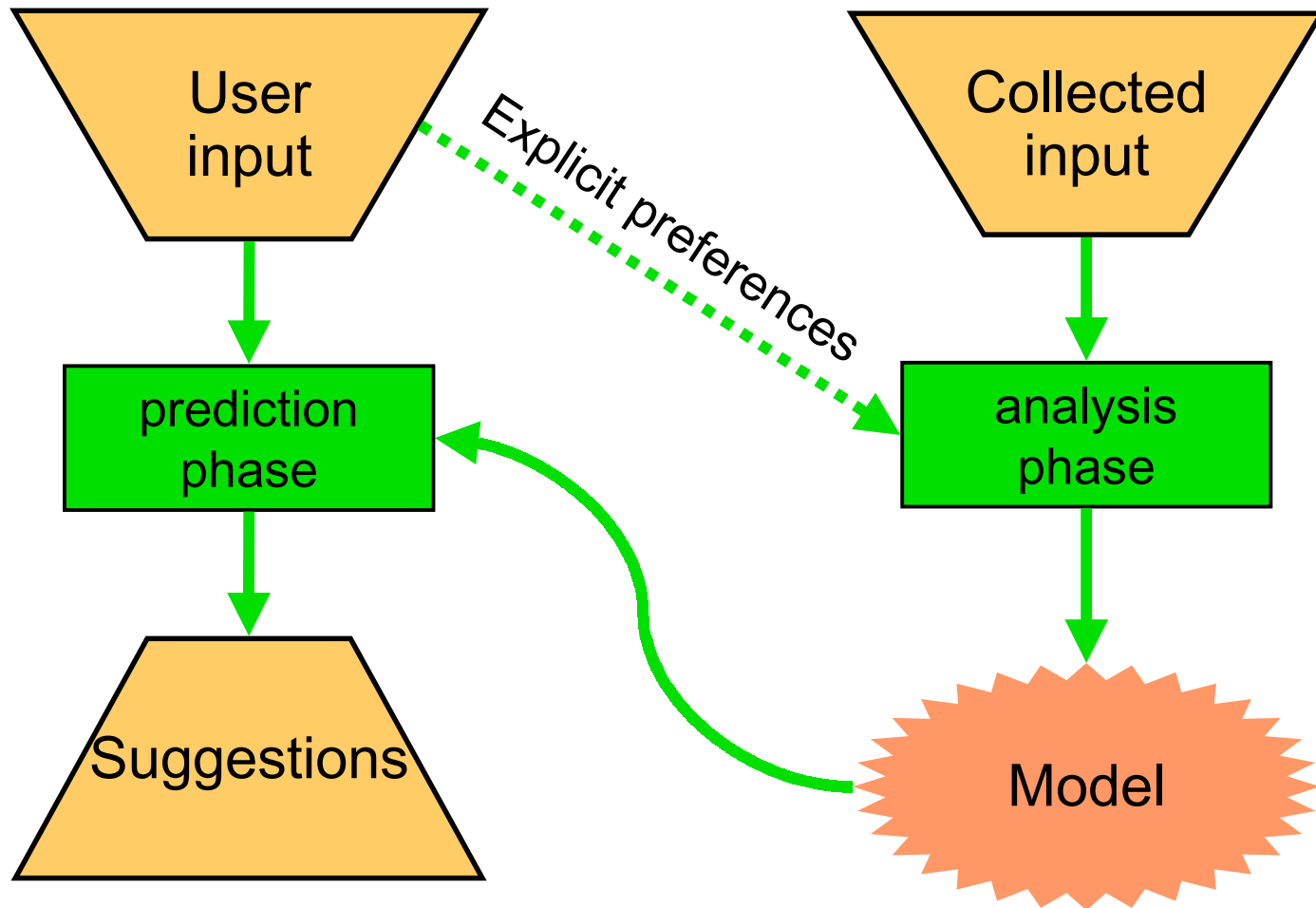
General-purpose search engines:

Hierarchical directories

Specialized search engines:

Search-by-example

- 1 Collaborative filtering
- 1 Meta-information



Google™ *Lots of projects*

- 1 Collaborative filtering seldom used in classic IR, big revival on the Web. Projects:
 - **PHOAKS** -- ATT labs → Web pages recommendation based on Usenet postings
 - **GAB** -- Bellcore → Web browsing
 - **Grouplens** -- U. Minnesota → Usenet newsgroups
 - **EachToEach** -- Compaq SRC → rating movies
 - ...

See <http://sims.berkeley.edu/resources/collab/>

1 The ranking schemes that we discussed are also a form of collaborative ranking!

- Connectivity = people vote with their links
- Usage = people vote with their clicks



1 These schemes are used only for a **global** model building. Can it be combined with **per-user** data?

Ideas:

- Consider the graph induced by the user's bookmarks.
- Profile the user -- see www.globalbrain.net
- Deal with privacy concerns!

General-purpose search engines:

Hierarchical directories

Specialized search engines:

Search-by-example

Collaborative filtering

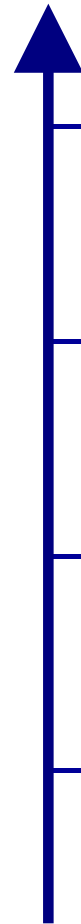
1 **Meta-information**

- Comparison of search engines
- Log statistics





**Difficulty of independent measurement;
Usefulness for Comparison**



Ideal measure: User satisfaction

Number of user requests

Quality of search engine index

Size of search engine index

1 Naïve Approaches

- Get a list of URLs from each search engine and compare
 - Not practical or reliable.
- Result Set Size Comparison
 - Reported sizes are approximate.
- ...

1 Better Approach

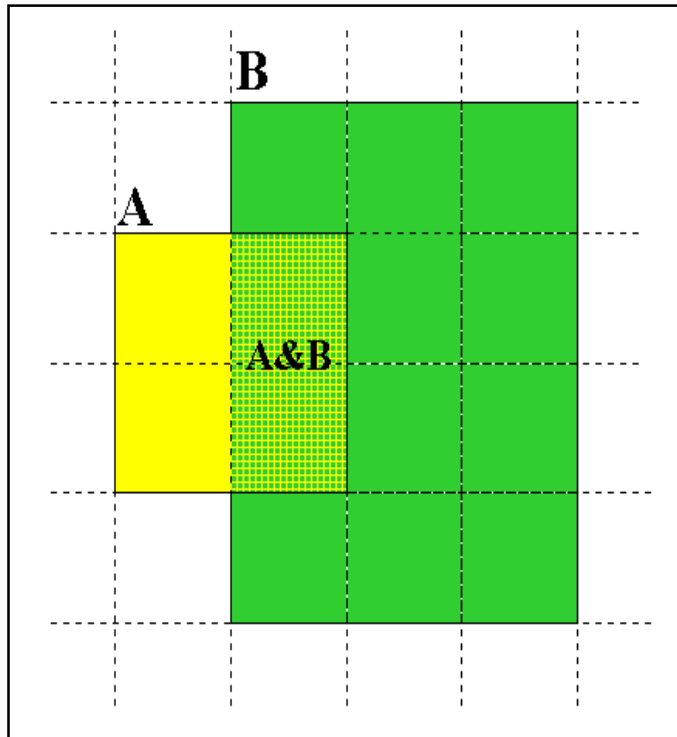
- Statistical Sampling

Google™ *URL sampling*

- 1 Ideal strategy: **Generate a random URL and check for containment in each index.**
- 1 Random URLs are hard to generate:
 - Random walks methods
 - Graph is directed
 - Stationary distribution is non-uniform
 - Must prove rapid mixing.
 - Pages in cache, query logs [LG'98a], etc.
 - Correlated to the interests of a particular group of users.
- 1 **A simple way: collect all pages on the Web and pick one at random.**



- 1 Search engines have the best crawlers -- why not exploit them?
- 1 Method:
 - Sample from each engine in turn
 - Estimate the relative sizes of two search engines
 - Compute absolute sizes from a reference point



Select pages randomly from A (resp. B)

Check if page contained in B (resp. A)

$$|A \cap B| \approx (1/2) * |A|$$

$$|A \cap B| \approx (1/6) * |B|$$

$$\therefore |B| \approx 3 * |A|$$

Two steps: (i) **Selecting** (ii) **Checking**

- 1 Generate random query
 - Build a lexicon of words that occur on the Web
 - Combine random words from lexicon to form queries
- 1 Get the first 100 query results from engine A
- 1 Select a random page out of the set
- 1 Distribution is biased -- the conjecture is that

$$\frac{\sum_{D \in A \cap B} p(D)}{\sum_{D \in A} p(D)} \sim \frac{|A \cap B|}{|A|}$$

where $p(D)$ is the probability that D is picked by this scheme

- 1 Create a “unique query” for the page:
 - Use 8 rare words.
 - E.g., for the Digital Systems Research Center Home Page:

[People Search](#) [Business Search](#)

Search the Web for documents in any language

+commercializing +lytton +nsl +crl +rad +accomplishments +mature -

search
refine

[Help](#) . [Preferences](#) . [New Search](#) . [Advanced Search](#)

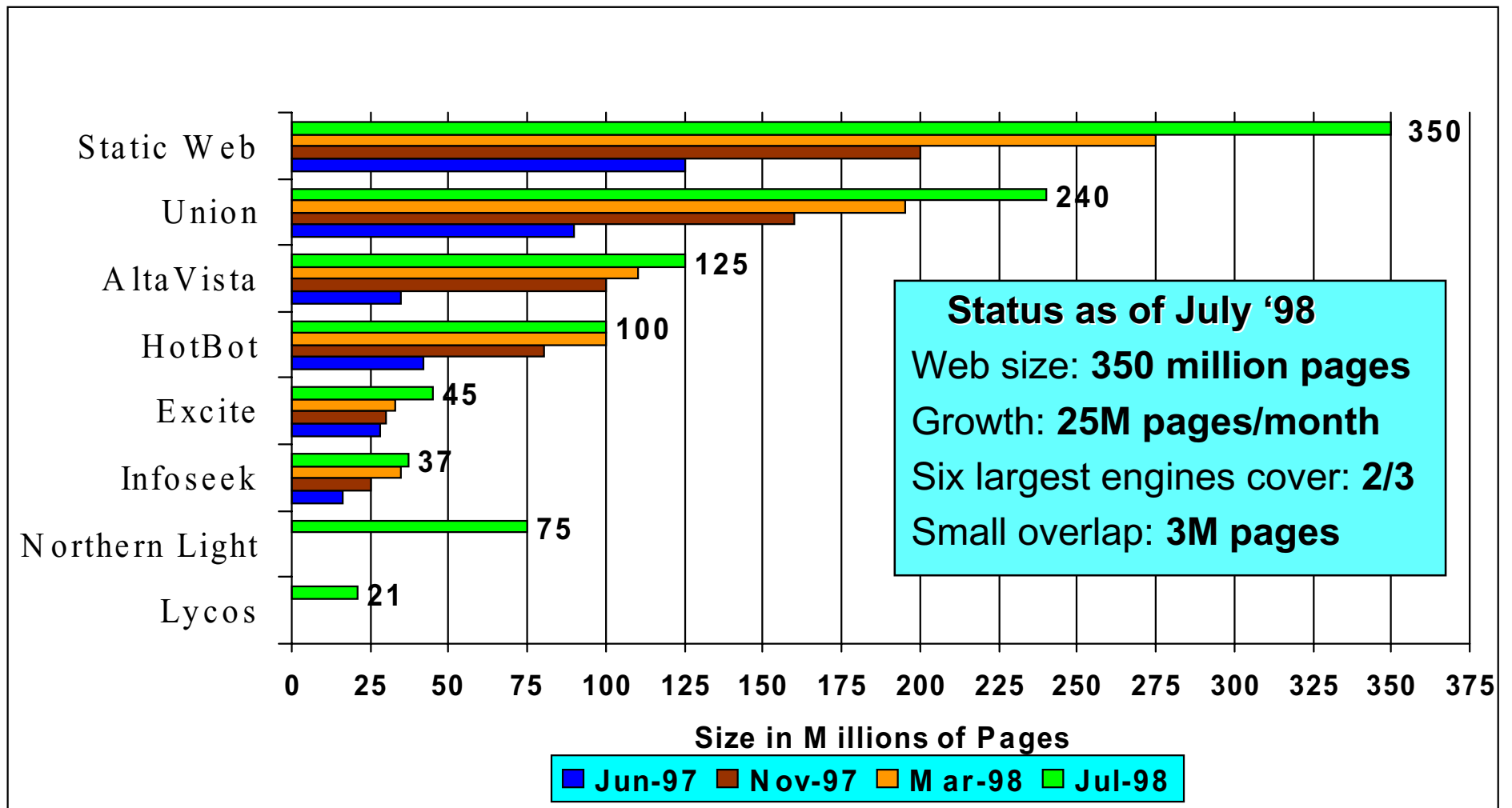
Click to find related books at [Amazon.com](#).

1 documents match your query.

1. [Systems Research Center - Home Page](#)
 The Systems Research Center (SRC) is one of four computer science research laboratories within Digital's Research and Advanced Development (RAD) group....
☞ <http://www.research.digital.com/src/home.html> - size 4K - 3-Oct-97 - English - [Translate](#)



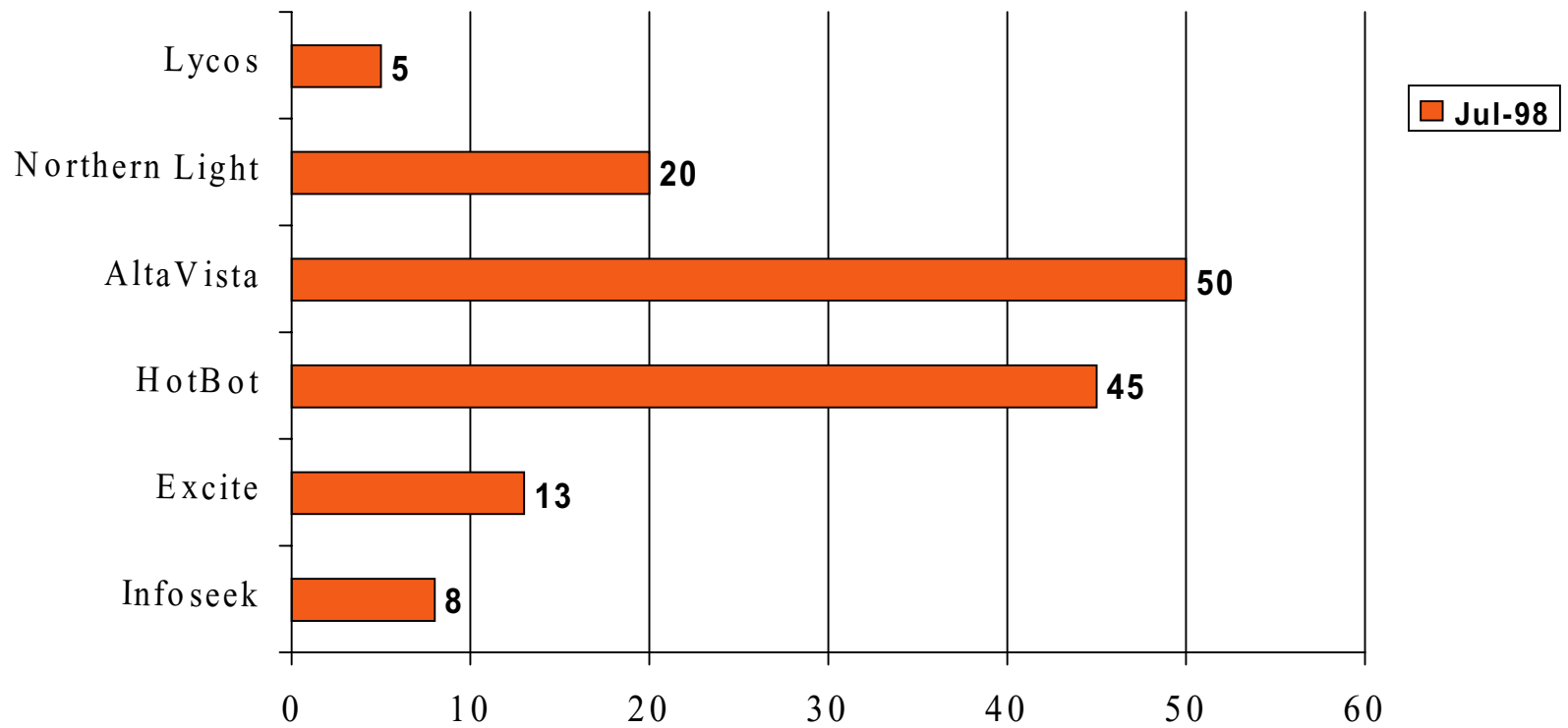
Results of the BB'98 study





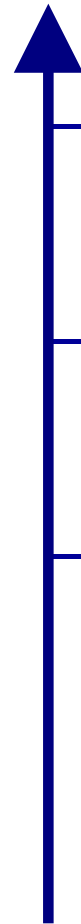
Crawling strategies are different!

Exclusive listings in millions of pages





**Difficulty of independent measurement;
Usefulness for Comparison**



Ideal measure: User satisfaction

Number of user requests

Quality of search engine index

Size of search engine index

Quality: A general definition

[HHMN'99]

- 1 Assign each page p a **weight** $w(p)$ such that

$$\sum_{\text{all } p} w(p) = 1$$

- 1 Can be thought of as probability distribution on pages

- 1 **Quality of a search engine index S is** $w(S) = \sum_{p \in S} w(p)$

- 1 **Example:**

– If w is same for all pages, weight is proportional to total size (in pages).

- 1 **Average page quality in index S is** $w(S)/|S|$.

- 1 We use: **weight** $w(p)$ of a page p is its **PageRank**

Suppose we can choose random pages according to w (so that page p appears with probability $w(p)$)

- 1 Choose a sample of pages $p_1, p_2, p_3 \dots p_n$
- 1 Check if the pages are in search engine index S
- 1 **Estimate for quality of index S** is the percentage of sampled pages that are in S , i.e.

$$\bar{w}(S) = \frac{1}{n} \sum_j I[p_j \in S]$$

where $I[p_j \text{ in } S] = 1$ if p_j is in S and 0 otherwise

Google™ *Missing pieces*

- 1 Sample pages according to the PageRank distribution.
- 1 Test whether page p is in search engine index S .
 - same methodology as [BB'98]



Sampling pages (almost) according to PageRank

- 1 Perform a random walk and select n random pages from it.
- 1 Problems:
 - Starting state bias: finite walk only approximates PageRank.
 - Can't jump to a random page; instead, jump to a random page on a random host seen so far.
- ⌘ Sampling pages according to a distribution that behaves similarly to PageRank, but it not identical to PageRank

Google *Experiments*

- 1 Performed two long random walks with $d=1/7$ starting at www.yahoo.com

	Walk 1	Walk2
length	18 hours	54 hours
attempted downloads	2,867,466	6,219,704
HTML pages successfully downloaded	1,393,265	2,940,794
unique HTML pages	509,279	1,002,745
sampled pages	1,025	1,100

- 1 Pages (or hosts) that are “highly-reachable” are visited often by the random walks
- 1 Initial bias for `www.yahoo.com` is reduced in longer walk
- 1 Results are consistent over the 2 walks
- 1 The average indegree of pages with indegree ≤ 1000 is high:
 - 53 in walk 1
 - 60 in walk 2



Most frequently visited pages

Page	Freq. Walk2	Freq. Walk1	Rank Walk1
www.microsoft.com/	3172	1600	1
www.microsoft.com/windows/ie/default.htm	2064	1045	3
www.netscape.com/	1991	876	6
www.microsoft.com/ie/	1982	1017	4
www.microsoft.com/windows/ie/download/	1915	943	5
www.microsoft.com/windows/ie/download/all.htm	1696	830	7
www.adobe.com/prodindex/acrobat/readstep.htm	1634	780	8
home.netscape.com/	1581	695	10
www.linkexchange.com/	1574	763	9
www.yahoo.com/	1527	1132	2

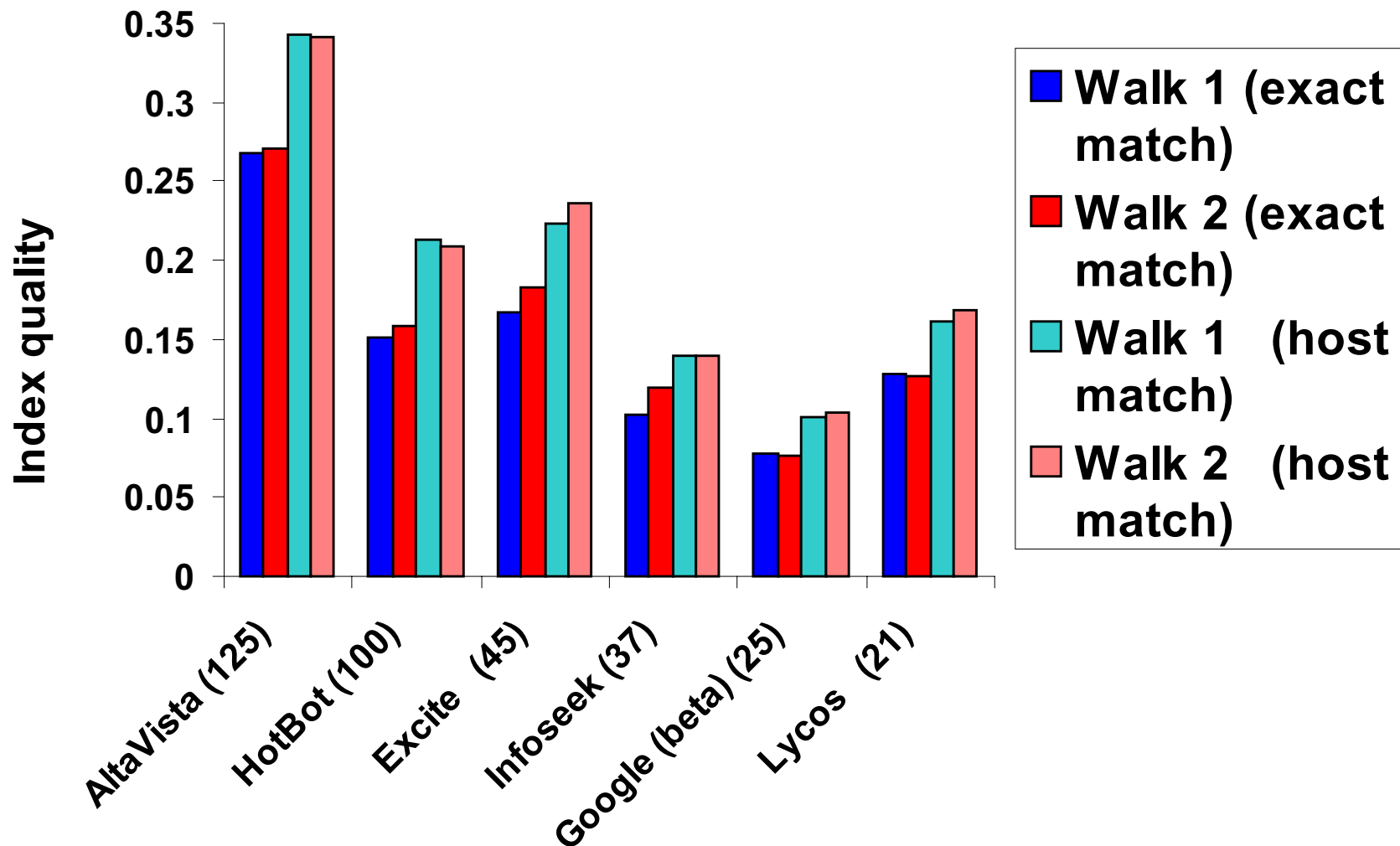


Most frequently visited hosts

Site	Frequency Walk 2	Frequency Walk 1	Rank Walk 1
www.microsoft.com	32452	16917	1
home.netscape.com	23329	11084	2
www.adobe.com	10884	5539	3
www.amazon.com	10146	5182	4
www.netscape.com	4862	2307	10
excite.netscape.com	4714	2372	9
www.real.com	4494	2777	5
www.lycos.com	4448	2645	6
www.zdnet.com	4038	2562	8
www.linkexchange.com	3738	1940	12
www.yahoo.com	3461	2595	7

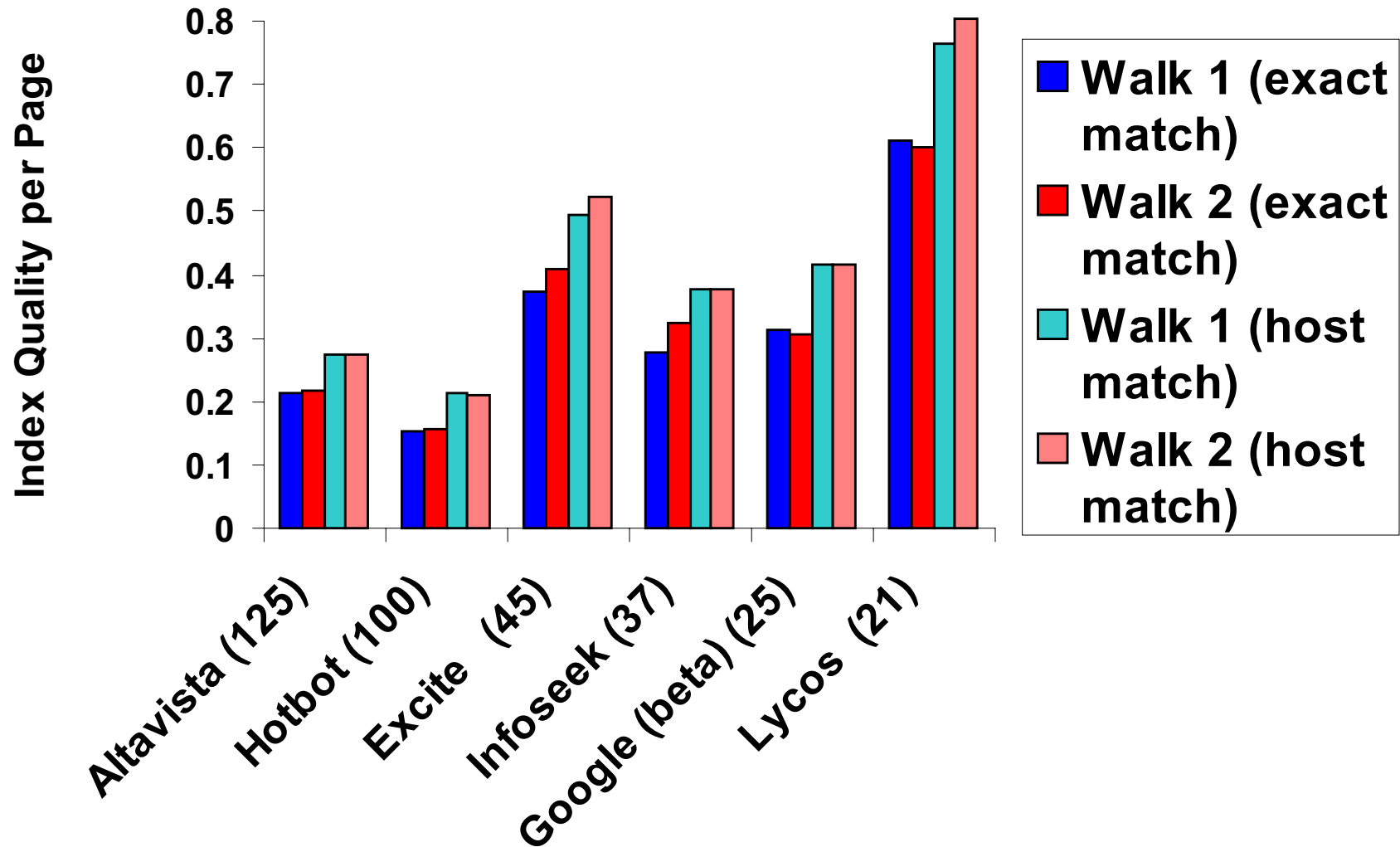


Results for index quality





Results for index quality/page



Google™ *Insights from the data*

- 1 Our approach appears consistent over repeated tests
- ⌘ Random walks are a useful tool
- 1 Quality is different from size for search engine indices
- ⌘ Some search engines are apparently trying to index
high quality pages

Google™ *Open problems*



- 1 Random page generation via random walks
- 1 Cryptography based approach: want **random pages from each engine but no cheating!** (page should be chosen u.a.r. from the actual index)
 - Each search engine can commit to the set of pages it has **without revealing it**
 - Need to ensure that this set is the same as the set actually indexed
 - Need efficient oblivious protocol to obtain random page from search engine
 - See [NP'98] for possible solution

General-purpose search engines:

Hierarchical directories

Specialized search engines:

Search-by-example

Collaborative filtering

1 **Meta-information**

Comparison of search engines

– Log statistics





How often do people view a page?

1 Problems:

- Web caches interfere with click counting
- cheating pays (advertisers pay by the click)

1 Solutions:

- naïve: forces caches to re-fetch for every click.
 - Lots of traffic, annoyed Web users
- extend HTML with counters [ML'97]
 - requires compliance, down caches falsify results.
- use sampling [P'97]
 - force refetches on random days
 - force refetches for random users and IP addresses
- cryptographic audit bureaus [NP'98a]

1 Commercial providers: 100hot, Media Matrix, Relevant Knowledge, ...

request = new query or new result screen of old query

session = a series of requests by one user close together in time

1 analyzed ~1B AltaVista requests consisting of:

- ~840 M non-empty requests
- ~575 M non-empty queries
 - ⇒ 1.5 requests per query in the average
- ~153 M unique non-empty queries
 - ⇒ query is repeated 3.8 times in the average, but 64% of queries occur only once
- ~285 M user sessions
 - ⇒ 2.9 requests and 2.0 queries per session in the average



Lots of things we didn't even touch...

- 1 **Clustering** = group similar items (documents or queries) together ↔ unsupervised learning
- 1 **Categorization** = assign items to predefined categories ↔ supervised learning
- 1 **Classic IR issues** that are not substantially different in the Web context:
 - Latent semantic indexing -- associate “concepts” to queries and documents and match on concepts
 - Summarization: abstract the most important parts of text content. (See [TS'98] for the Web context)
- 1 **Many others ...**

- 1 We talked mostly about IR methods and tools that
 - take advantage of the Web particularities
 - mitigate some of the difficulties
- 1 Web IR offers plenty of interesting problems...
 - ... but not on a silver platter
- 1 Almost every area of algorithms research is relevant
- 1 Great need for good algorithm engineering!



Acknowledgements

- 1 An earlier version of this talk was created in collaboration with Andrei Broder and was presented at the 39th IEEE Symposium on Foundations in Computer Science 1998.
- 1 Thanks for useful comments to Krishna Bharat, Moses Charikar, Edith Cohen, Jeff Dean, Uri Feige, Puneet Kumar, Hannes Marais, Michael Mitzenmacher, Prabhakar Raghavan, and Lyle Ramshaw.