

Web mapping, modeling, mining, and mingling

Tutorial for the 2006 MITACS Winter School
Modelling and Mining of Networked Information Spaces

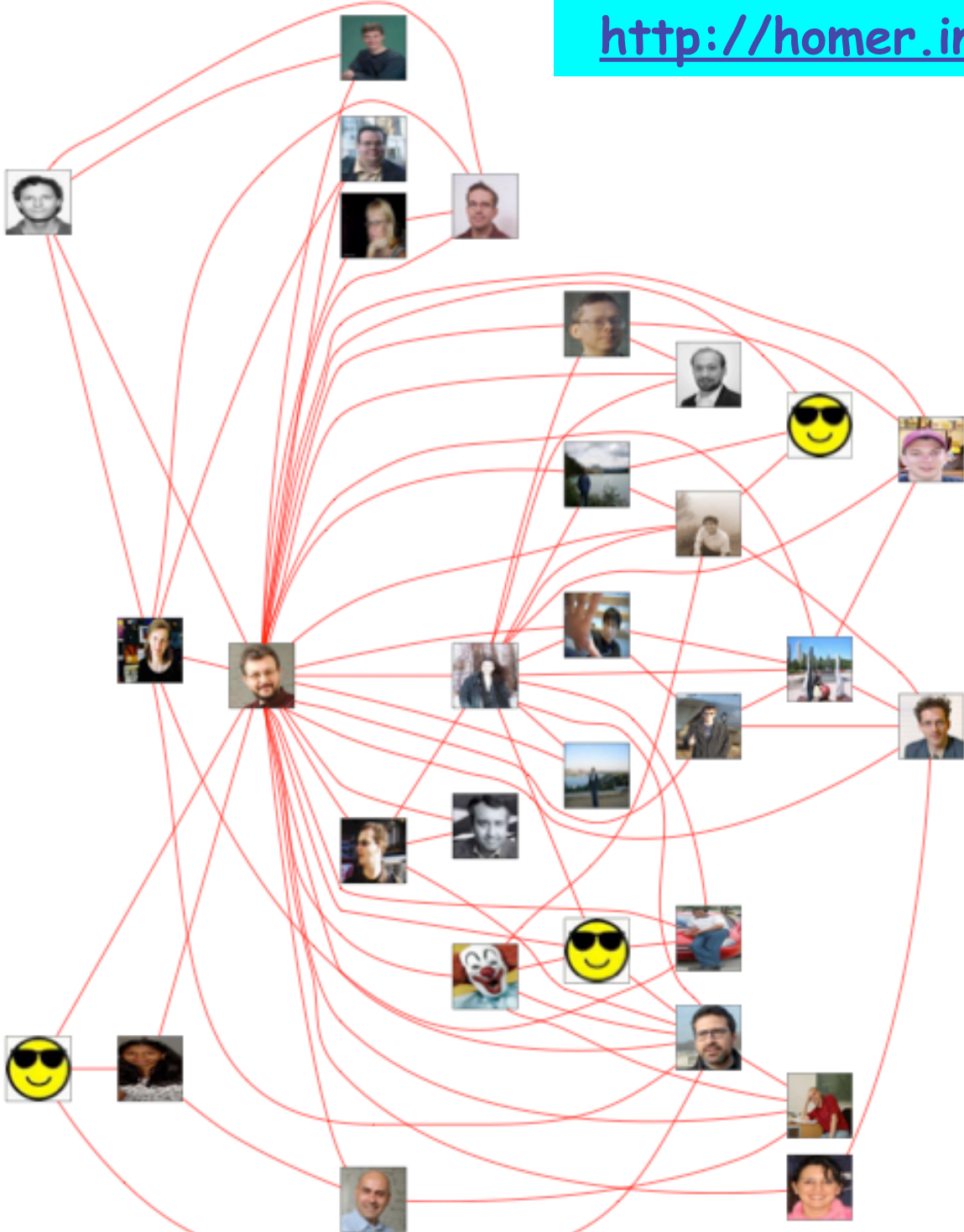
Filippo Menczer

Informatics & Computer Science
Indiana University, Bloomington

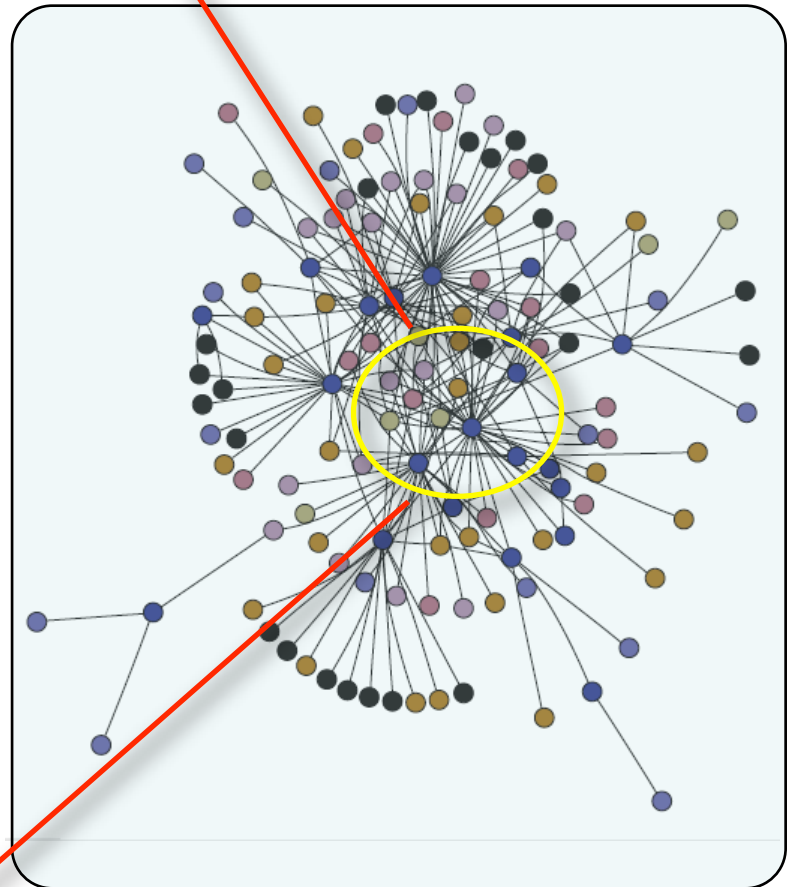
Indiana University Department of
ComputerScience

Indiana University School of
informatics

<http://homer.informatics.indiana.edu/~nan/>



NaN:
Networks
& agents
Network

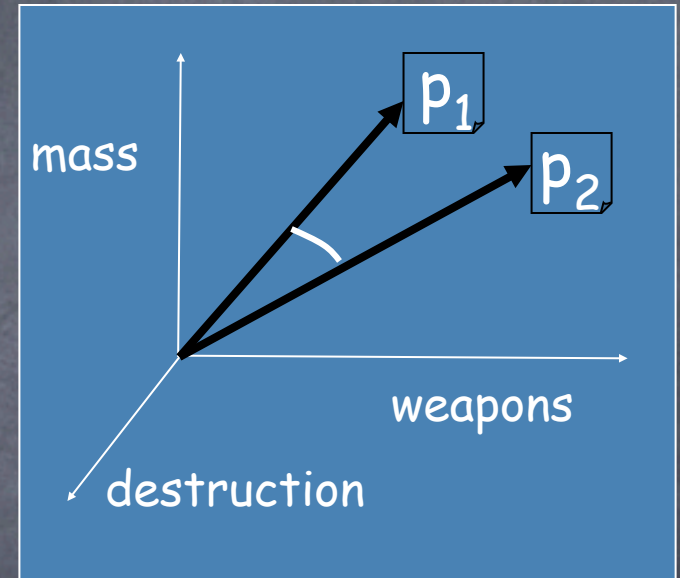


Outline

- Mapping
 - > Topical locality
 - > Content, link, and semantic topologies in the Web
- Modeling
 - > How the Web evolves and why content matters
 - > Consequences for navigation and crawling
- Mining
 - > Topical Web crawlers
 - > Adaptive, intelligent crawling techniques
- Mingling
 - > Social Web search & recommendation
 - > Distributed collaborative peer search

The Web as a text corpus

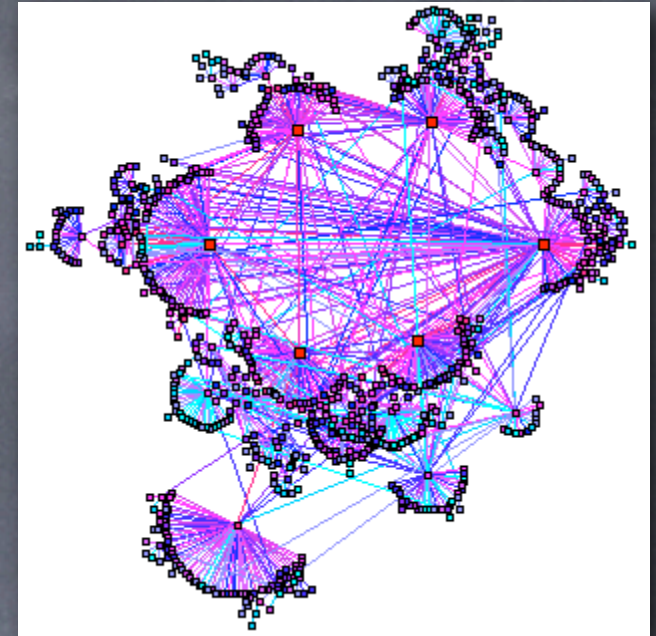
- Pages close in **word vector space** tend to be related
- Cluster hypothesis (van Rijsbergen 1979)
- The WebCrawler (Pinkerton 1994)
- The whole first generation of search engines



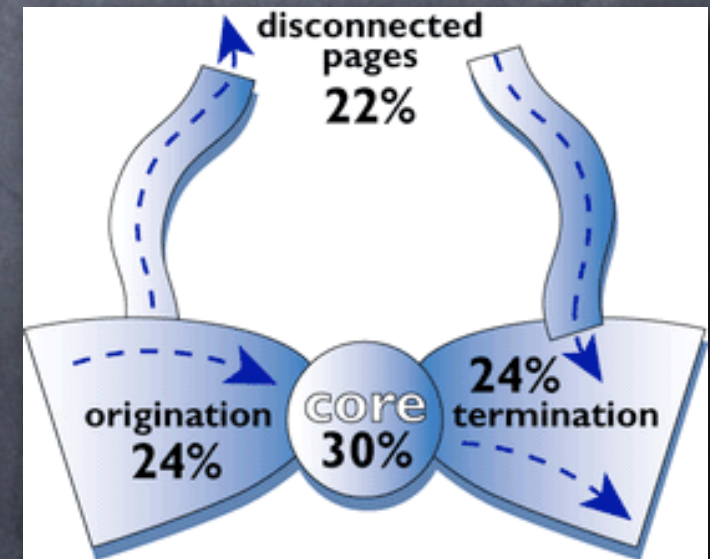
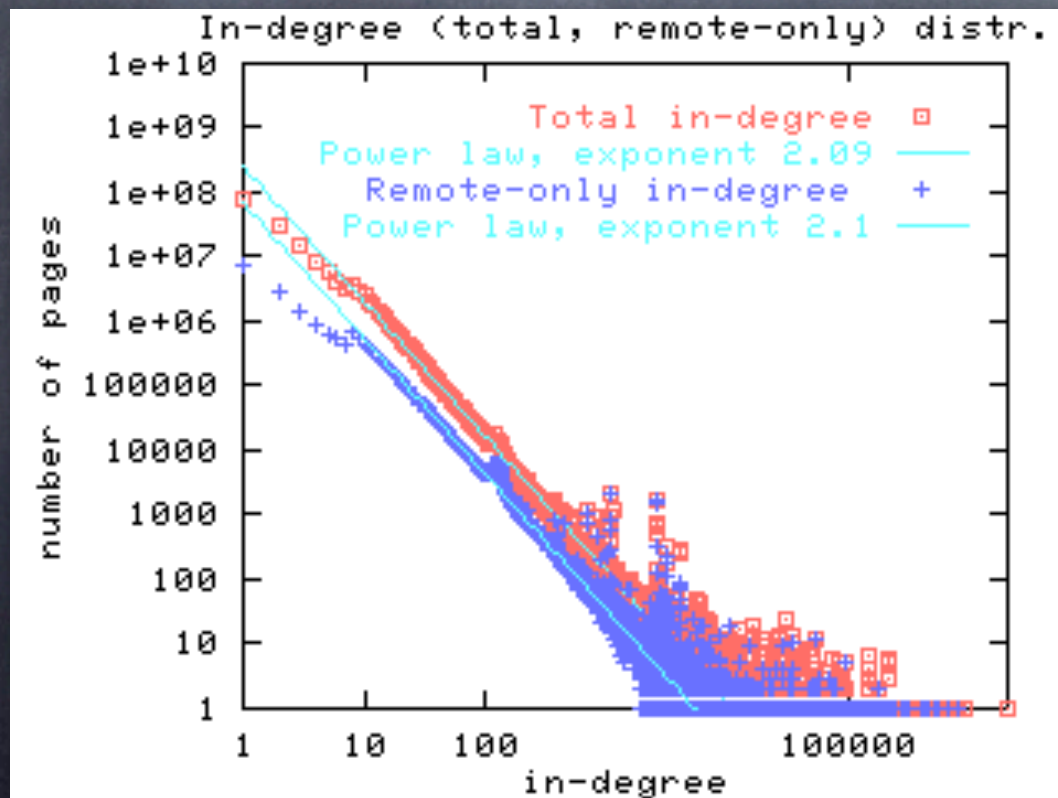
Enter the Web's link structure

$$p(i) = \frac{\alpha}{N} + (1 - \alpha) \sum_{j:j \rightarrow i} \frac{p(j)}{|\ell : j \rightarrow \ell|}$$

Brin & Page 1998

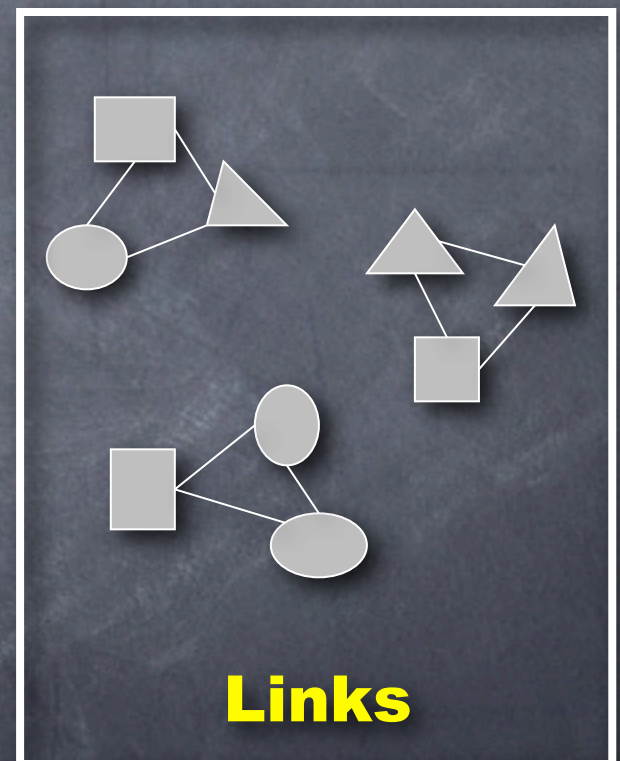
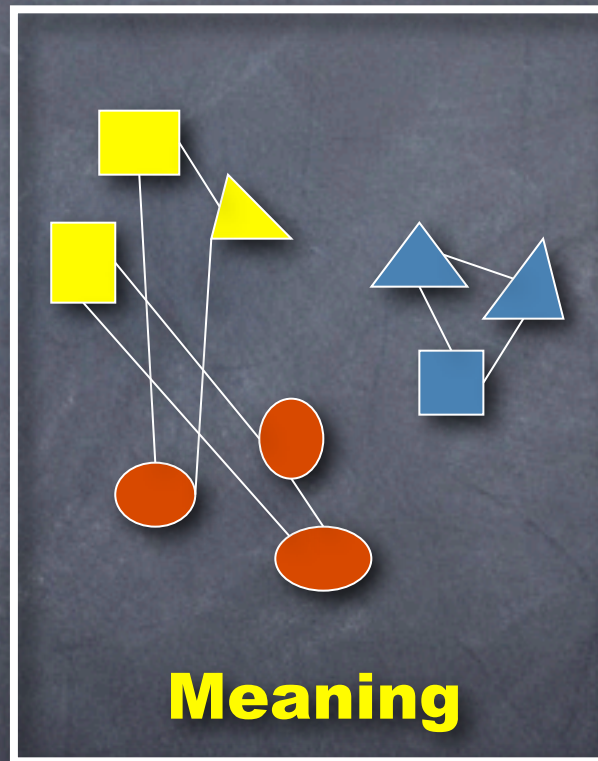
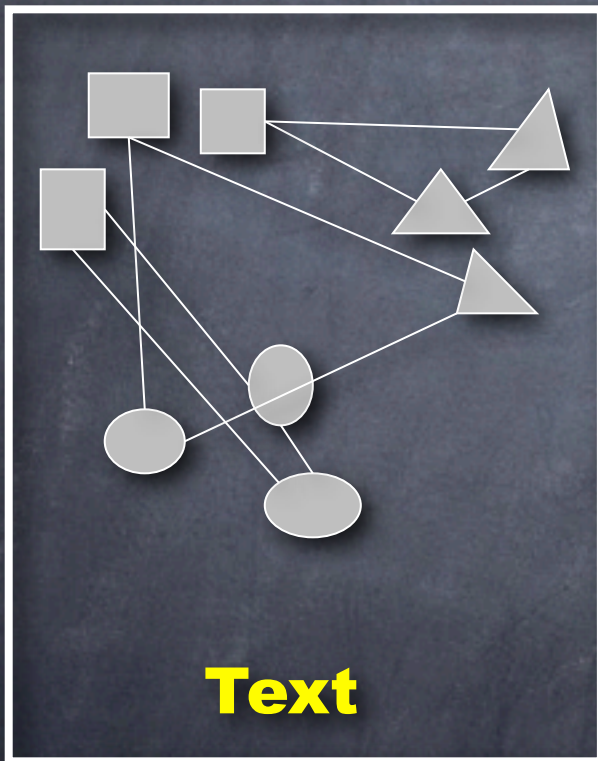


Barabasi & Albert 1999



Broder & al. 2000

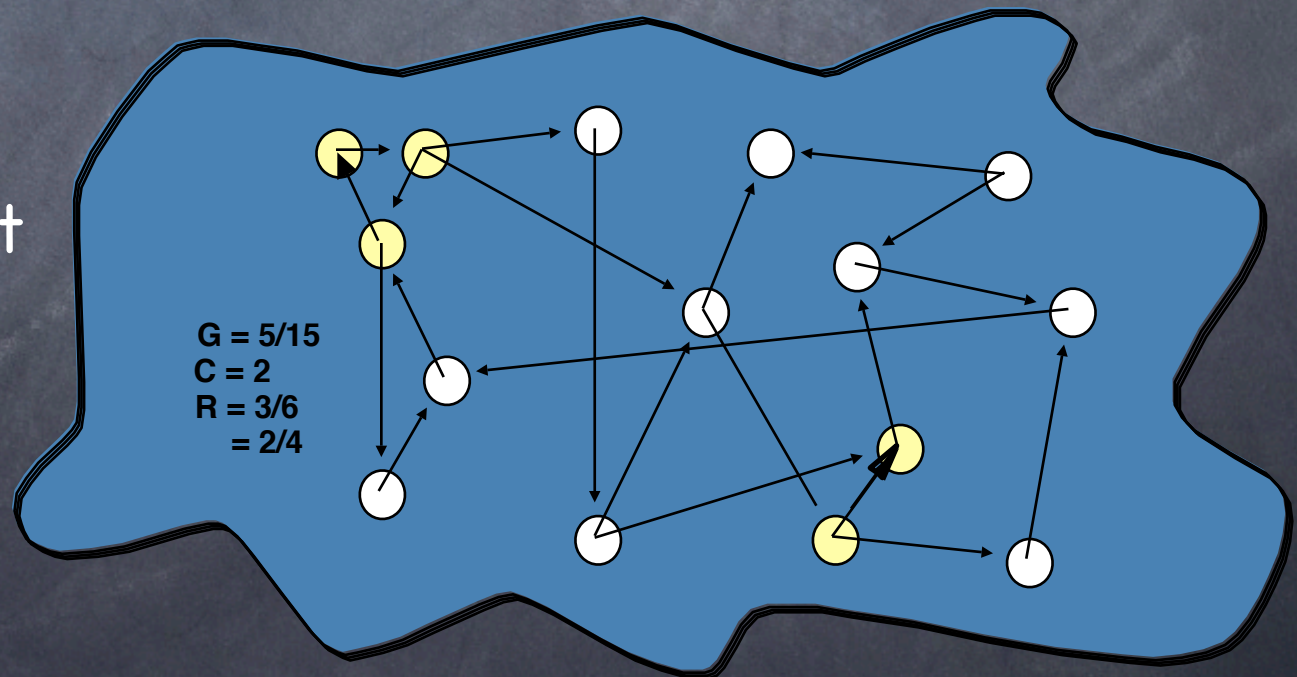
Three network topologies



The "link-cluster" conjecture

- Connection between **semantic** topology (topicality or relevance) and **link** topology (hypertext)
- $G = \Pr[\text{rel}(p)] \sim$ fraction of relevant pages (generality)
- $R = \Pr[\text{rel}(p) \mid \text{rel}(q) \text{ AND } \text{link}(q,p)]$
- Related nodes are "clustered" if **$R > G$** (modularity)

- Necessary and sufficient condition for a random crawler to find pages related to start points



Link-cluster conjecture

- Stationary hit rate for a random crawler:

$$\eta(t+1) = \eta(t) \cdot R + (1 - \eta(t)) \cdot G \geq \eta(t)$$

$$\eta \xrightarrow{t \rightarrow \infty} \eta^* = \frac{G}{1 - (R - G)}$$

$$\eta^* > G \Leftrightarrow R > G$$

$$\frac{\eta^*}{G} - 1 = \frac{R - G}{1 - (R - G)}$$

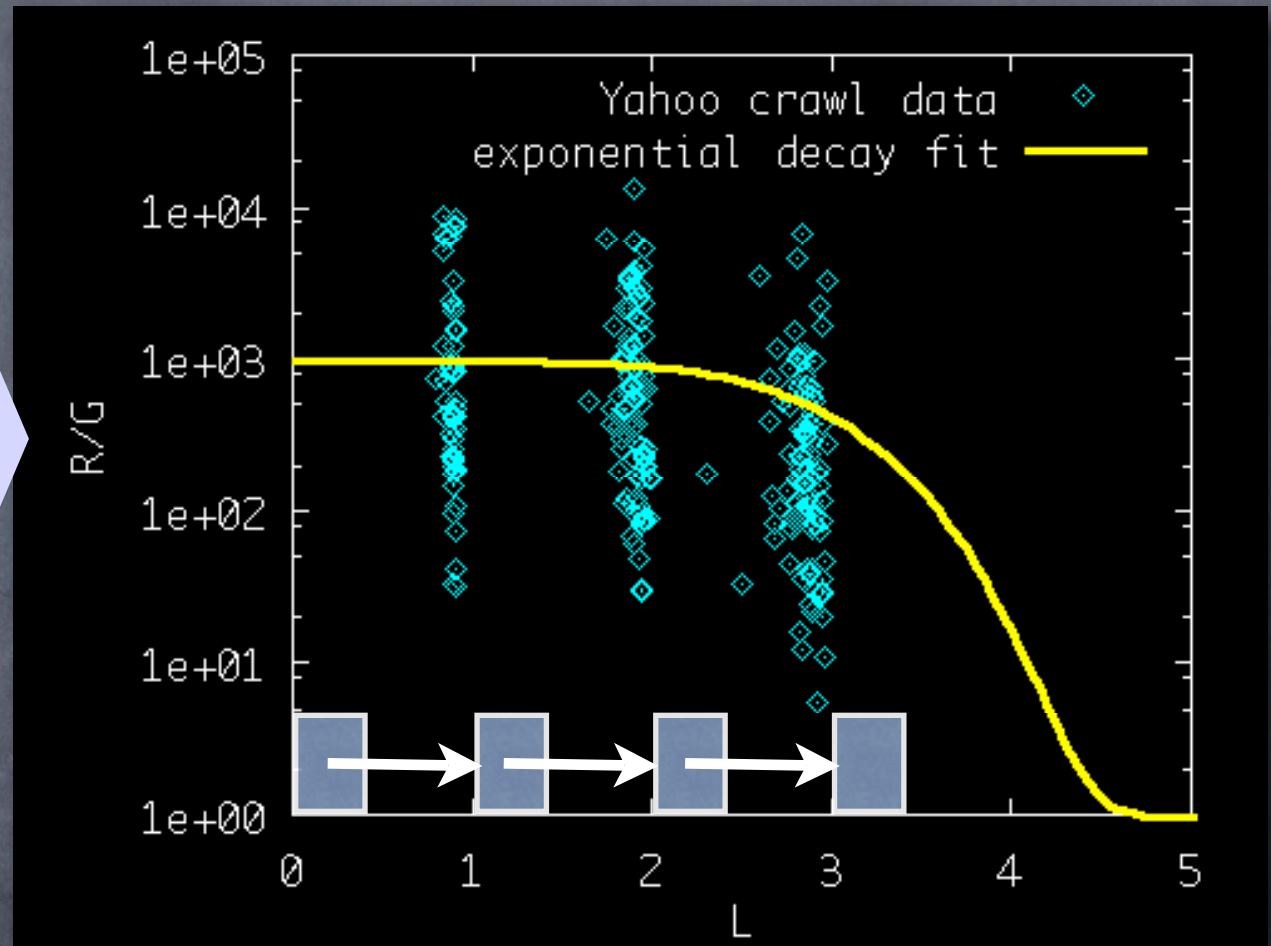
Conjecture

Value added

$$\frac{R(q, \delta)}{G(q)} \equiv \frac{\Pr[rel(p) \mid rel(q) \wedge \|path(q, p)\| \leq \delta]}{\Pr[rel(p)]}$$

Link-cluster conjecture

- Pages that **link to each other** tend to be related
- Preservation of **semantics** (meaning)
- A.k.a. **topic drift**

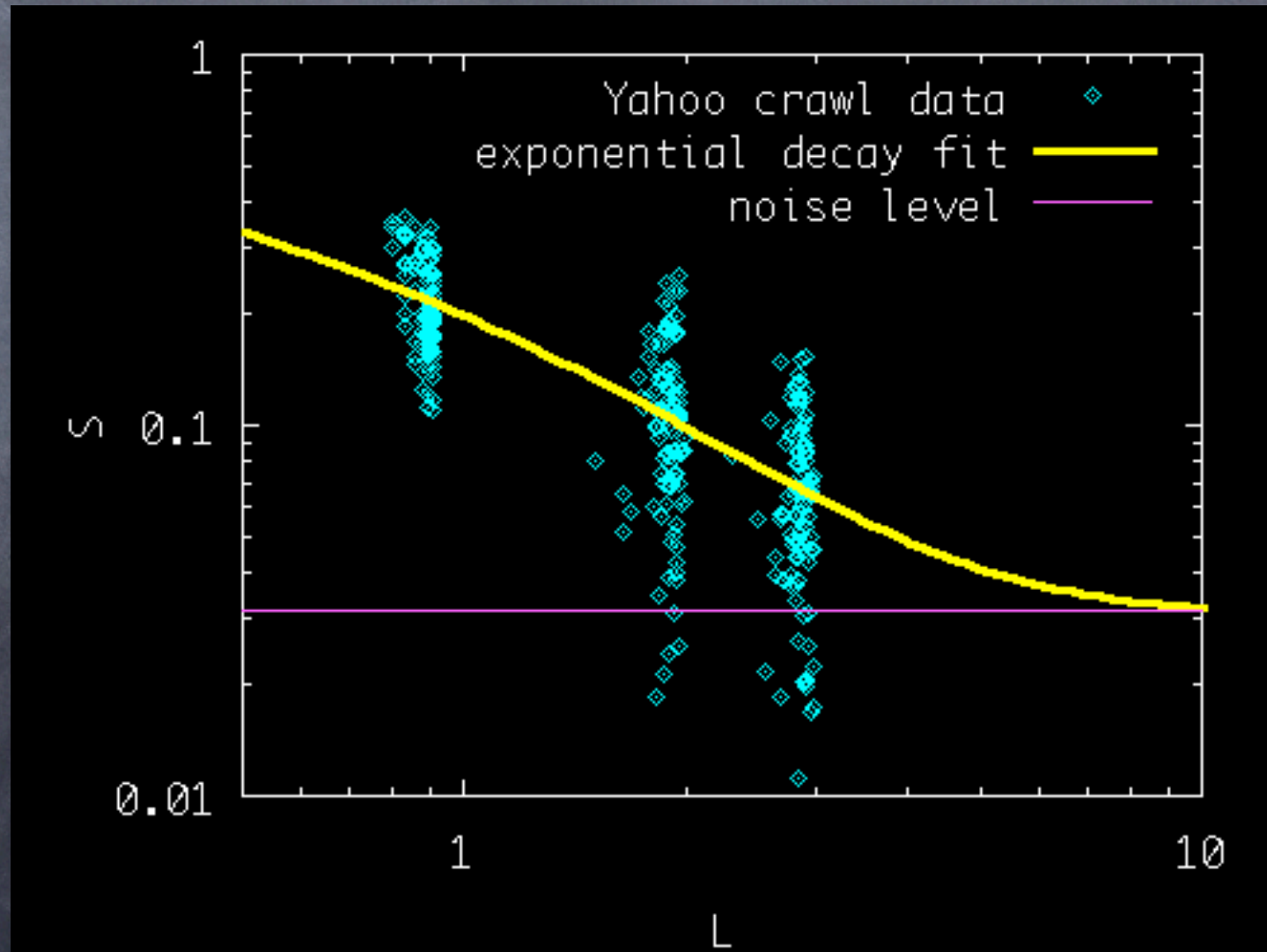


$$L(q, \delta) \equiv \frac{\sum_{\{p: \|path(q, p)\| \leq \delta\}} \|path(q, p)\|}{|\{p: \|path(q, p)\| \leq \delta\}|}$$

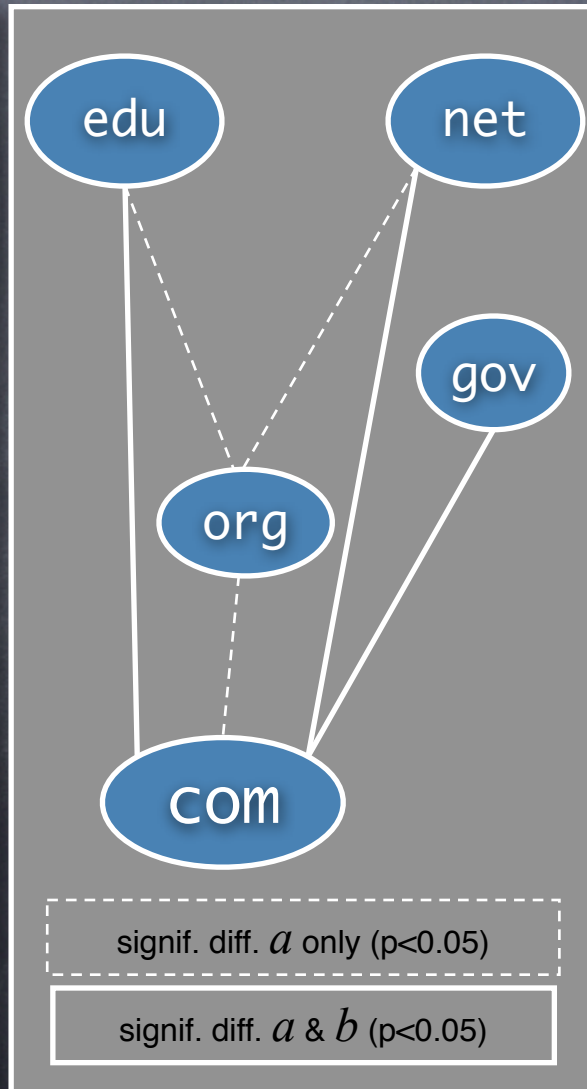
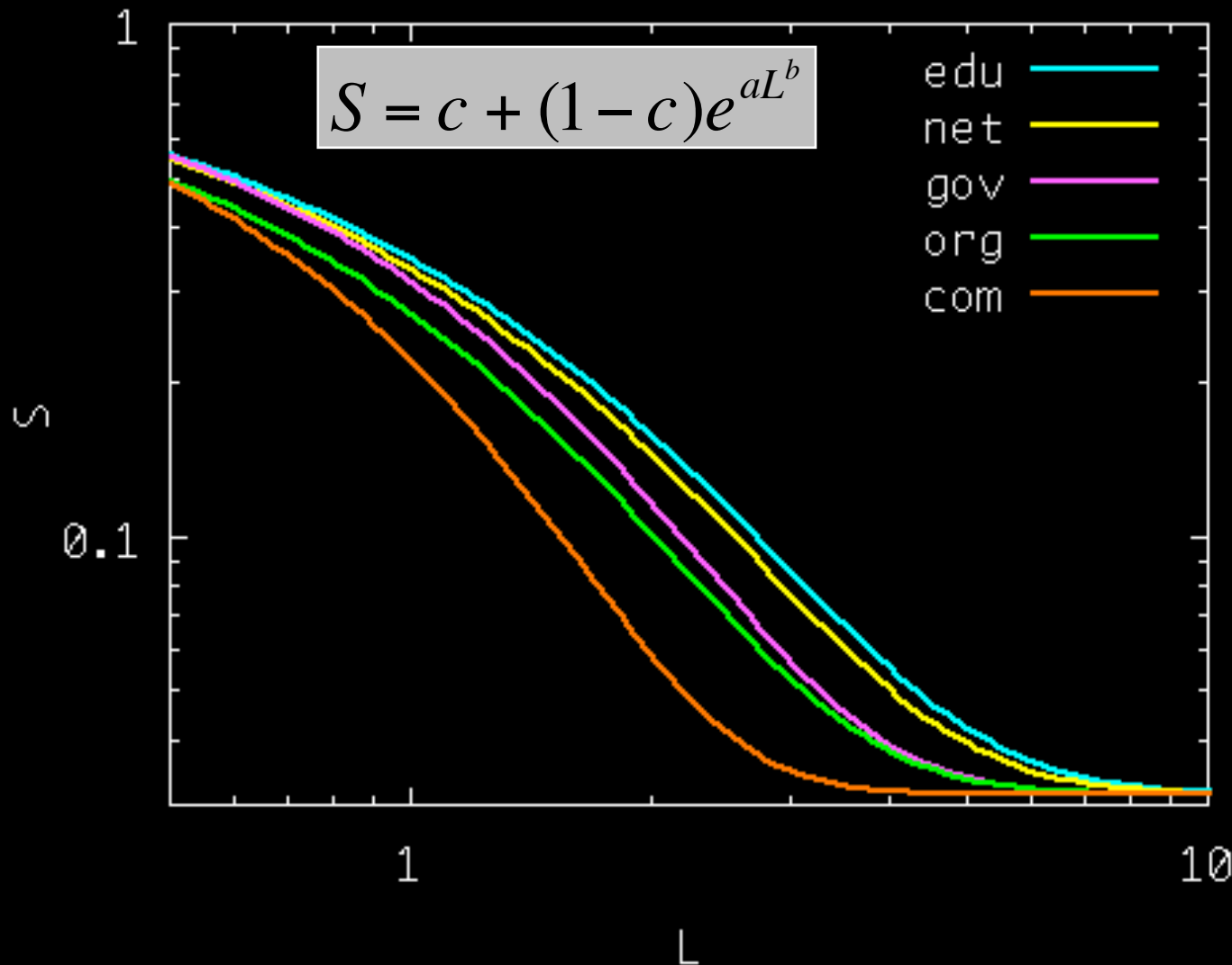
The “link-content” conjecture

- Correlation of **lexical** and **linkage** topology
- **$L(\delta)$** : average link distance
- **$S(\delta)$** : average similarity to start (topic) page from pages up to distance δ
- Correlation **$\rho(L, S) = -0.76$**

$$S(q, \delta) \equiv \frac{\sum_{\{p: \|path(q, p)\| \leq \delta\}} sim(q, p)}{|\{p: \|path(q, p)\| \leq \delta\}|}$$



Heterogeneity of link-content correlation



Discussion

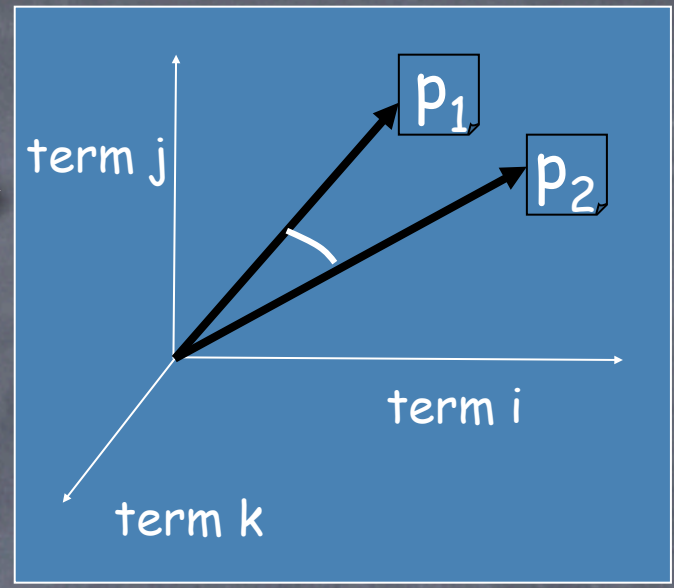
Topic drift:
Myth or reality?

Mapping the relationship between links, content, and **semantic** topologies

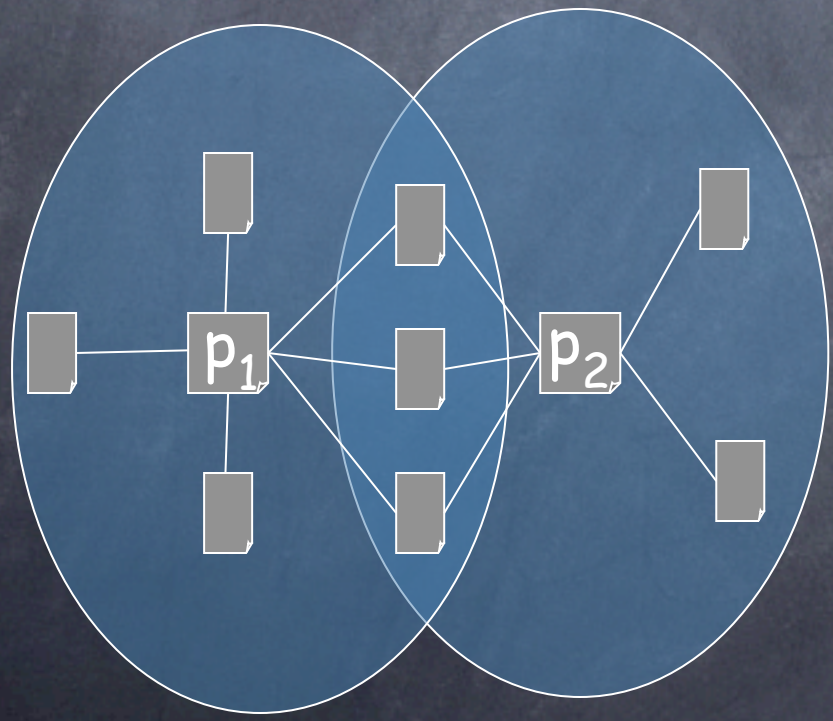
- Given any pair of pages, need 'similarity' or 'proximity' metric for each topology:
 - **Content**: textual/lexical (cosine) similarity
 - **Link**: co-citation/bibliographic coupling
 - **Semantic**: relatedness inferred from manual classification
- Data: Open Directory Project (**dmoz.org**)
 - ~ 1 M pages after cleanup
 - ~ 1.3×10^{12} page pairs!

$$\sigma_c(\vec{p}_1, \vec{p}_2) = \frac{\vec{p}_1 \cdot \vec{p}_2}{\|\vec{p}_1\| \cdot \|\vec{p}_2\|}$$

Content similarity

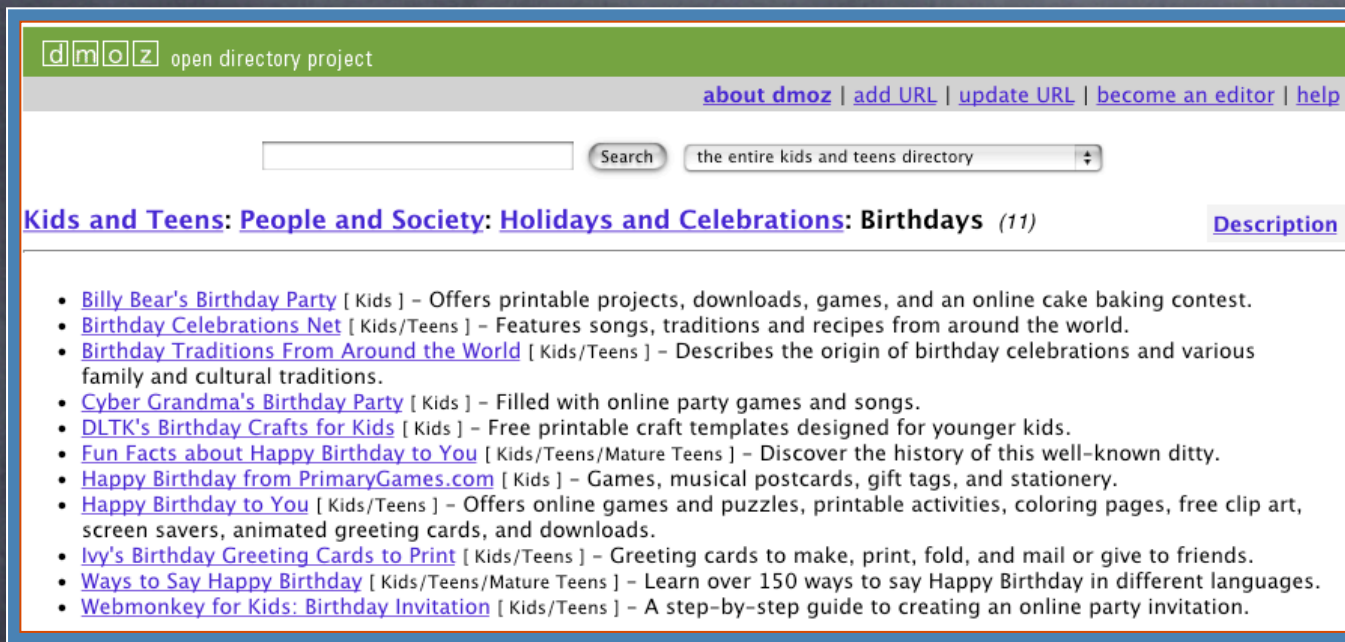


Link similarity

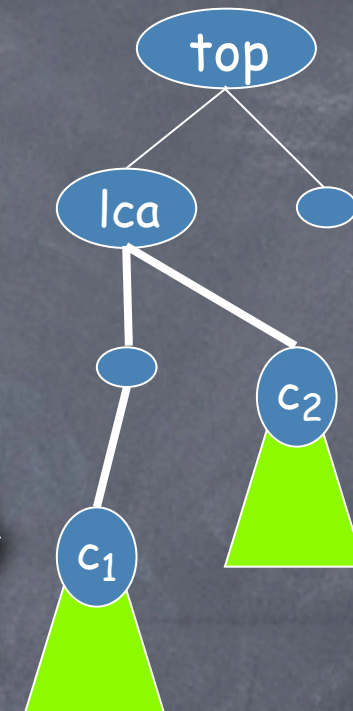


$$\sigma_l(p_1, p_2) = \frac{|U_{p_1} \cap U_{p_2}|}{|U_{p_1} \cup U_{p_2}|}$$

Semantic similarity



The screenshot shows the DMOZ directory page for "Kids and Teens: People and Society: Holidays and Celebrations: Birthdays". The page features a search bar, navigation links, and a list of 11 items related to birthdays, such as "Billy Bear's Birthday Party" and "Birthday Celebrations Net".

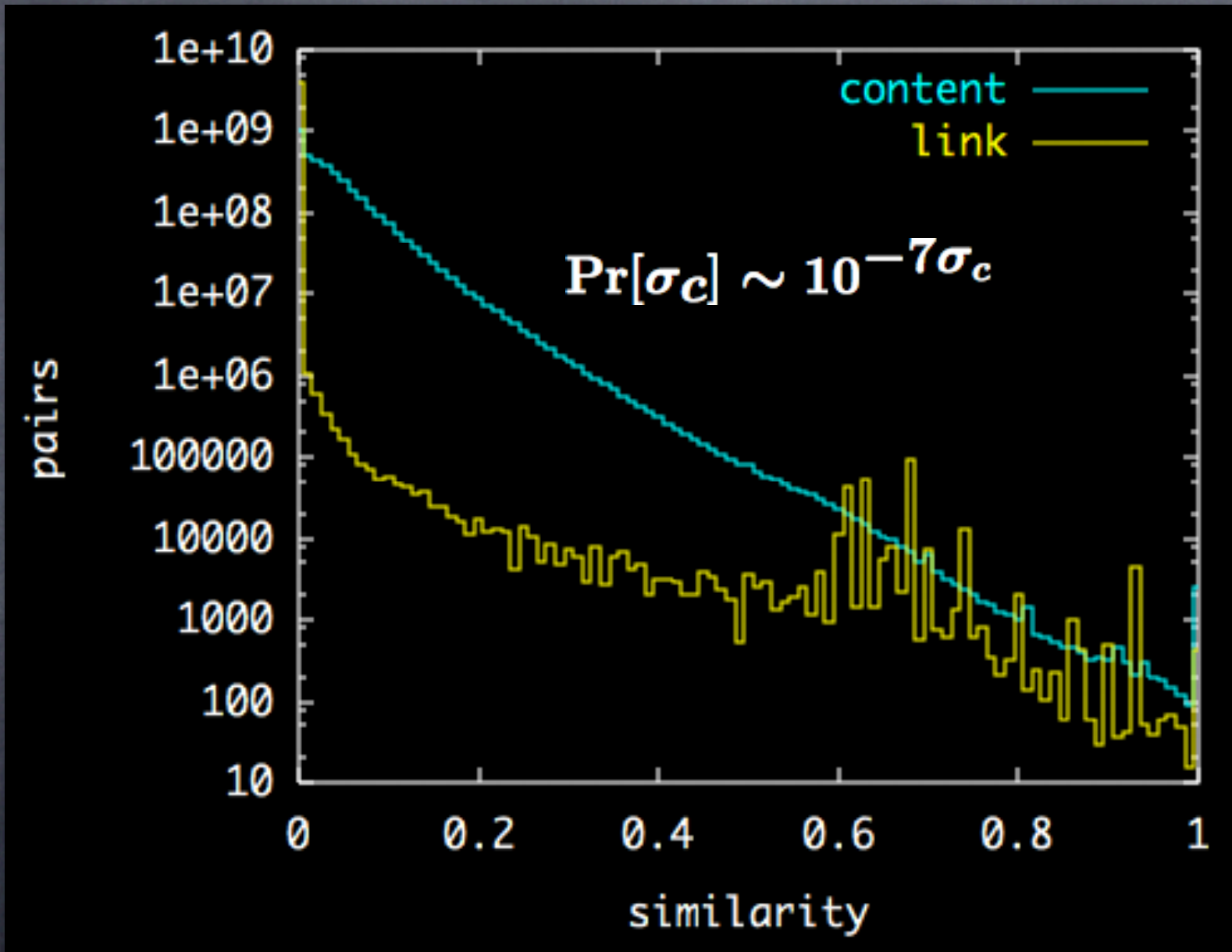


- Information-theoretic measure based on classification tree (Lin 1998)

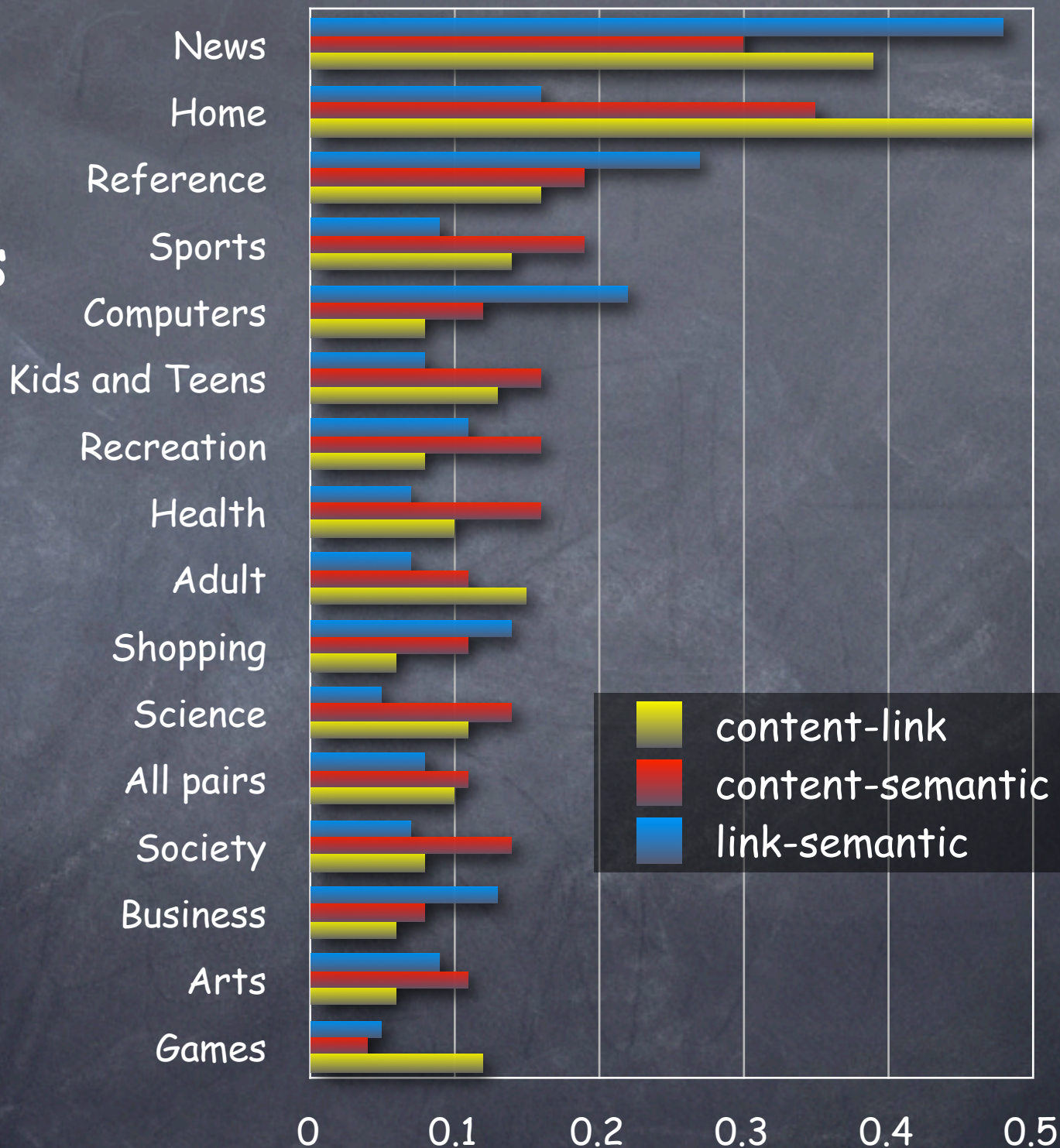
$$\sigma_s(c_1, c_2) = \frac{2 \log \Pr[lca(c_1, c_2)]}{\log \Pr[c_1] + \log \Pr[c_2]}$$

- Classic path distance in special case of balanced tree

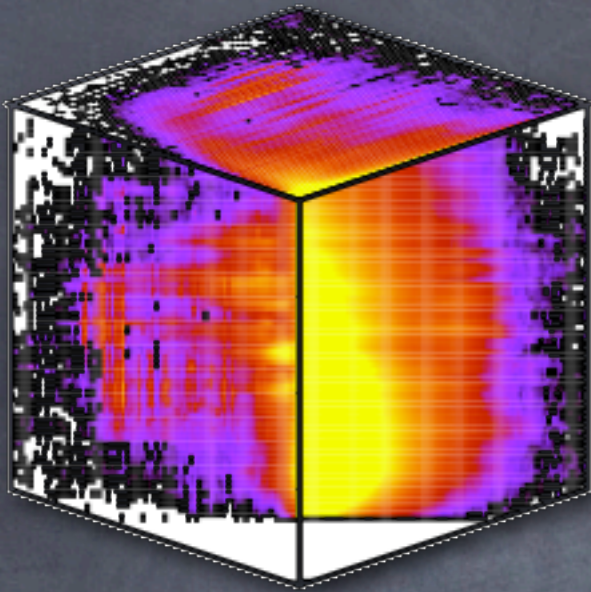
Individual metric distributions



Correlations between similarities



*IEEE Internet
Computing 2005*



$$\text{Precision} = \frac{|\text{Retrieved \& Relevant}|}{|\text{Retrieved}|}$$

$$\text{Recall} = \frac{|\text{Retrieved \& Relevant}|}{|\text{Relevant}|}$$

$$P(s_c, s_l) = \frac{\sum_{\{p, q: \sigma_c = s_c, \sigma_l = s_l\}} \sigma_s(p, q)}{|\{p, q: \sigma_c = s_c, \sigma_l = s_l\}|}$$

Averaging
semantic
similarity

$$R(s_c, s_l) = \frac{\sum_{\{p, q: \sigma_c = s_c, \sigma_l = s_l\}} \sigma_s(p, q)}{\sum_{\{p, q\}} \sigma_s(p, q)}$$

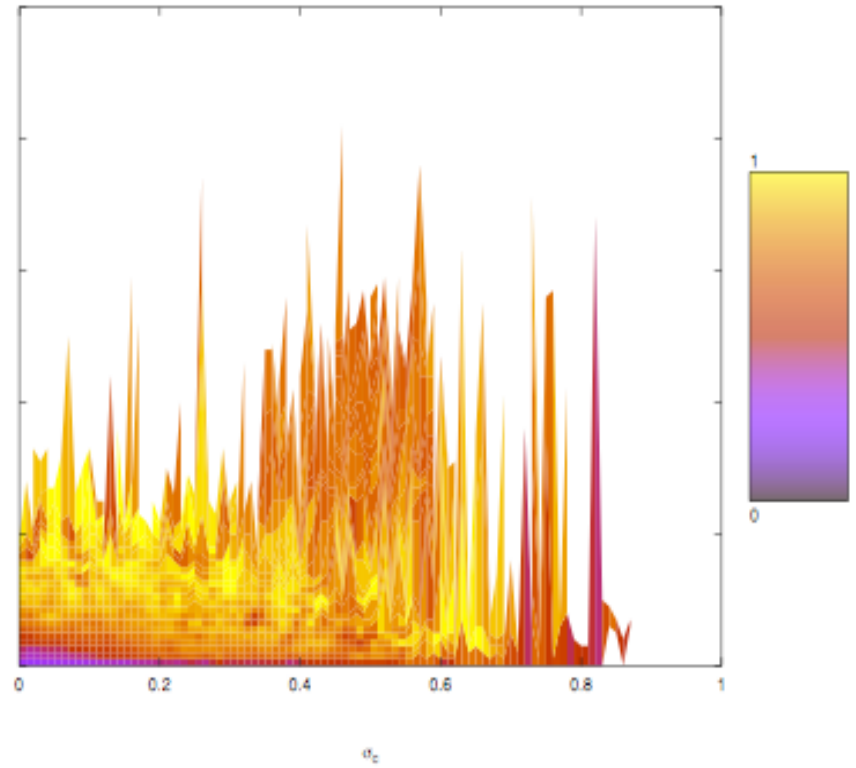
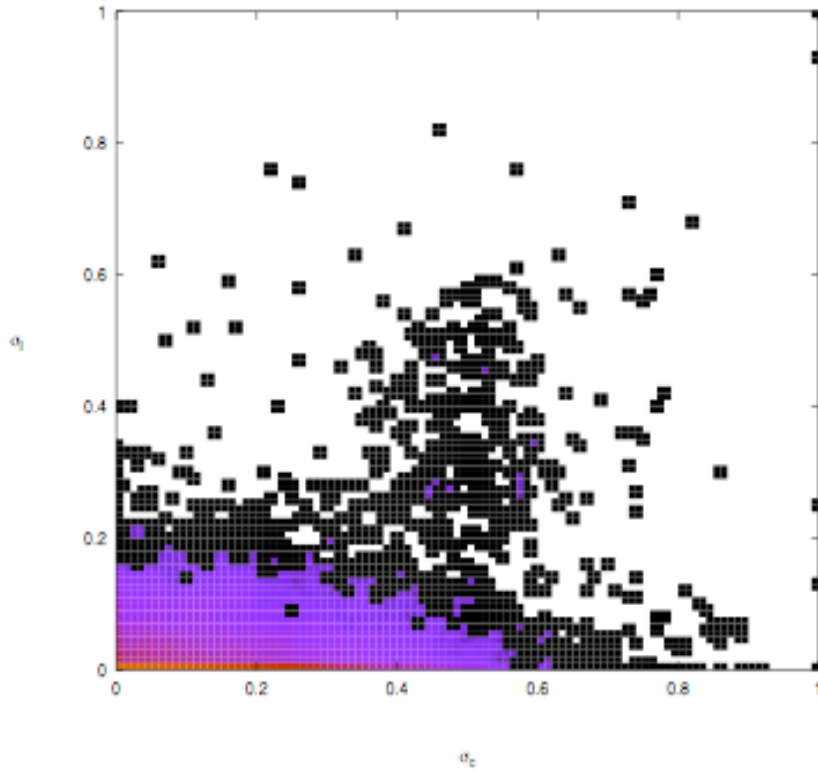
Summing
semantic
similarity

Business

σ_ℓ

log Recall

Precision



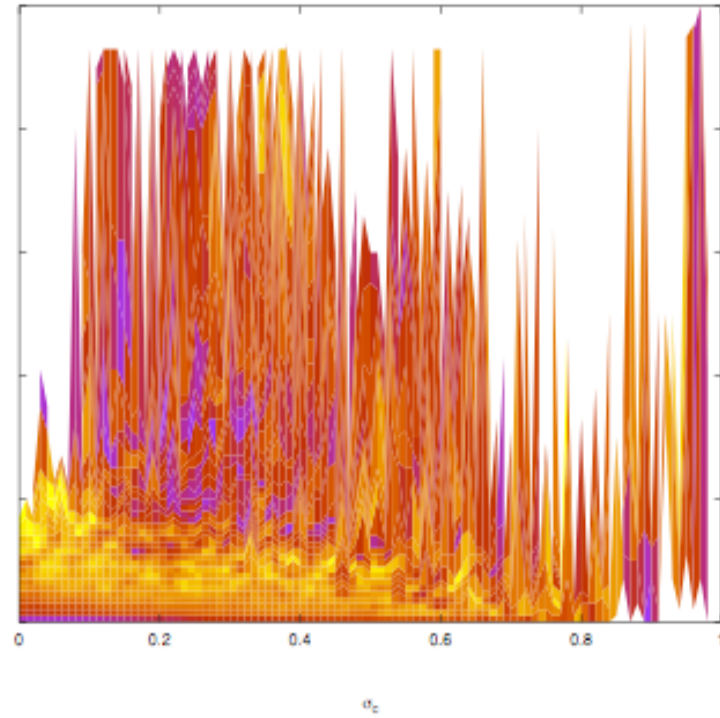
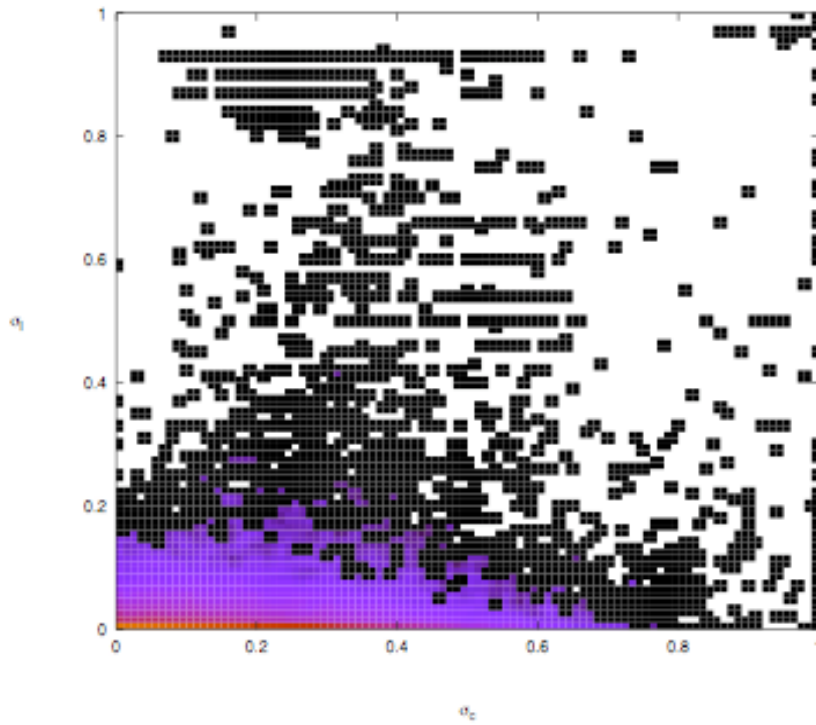
σ_c

Science

log Recall

Precision

σ_ℓ



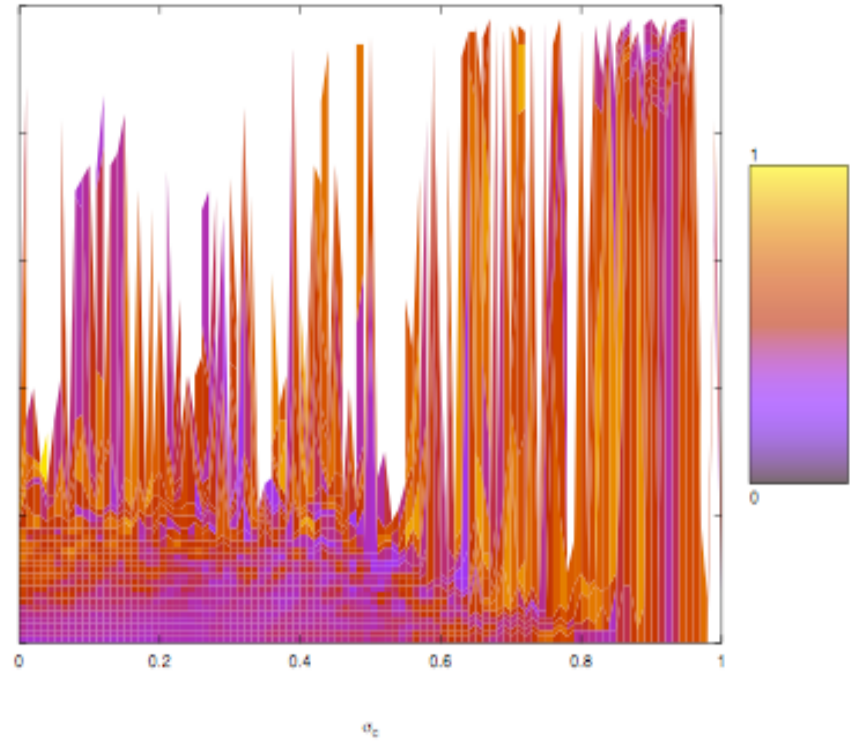
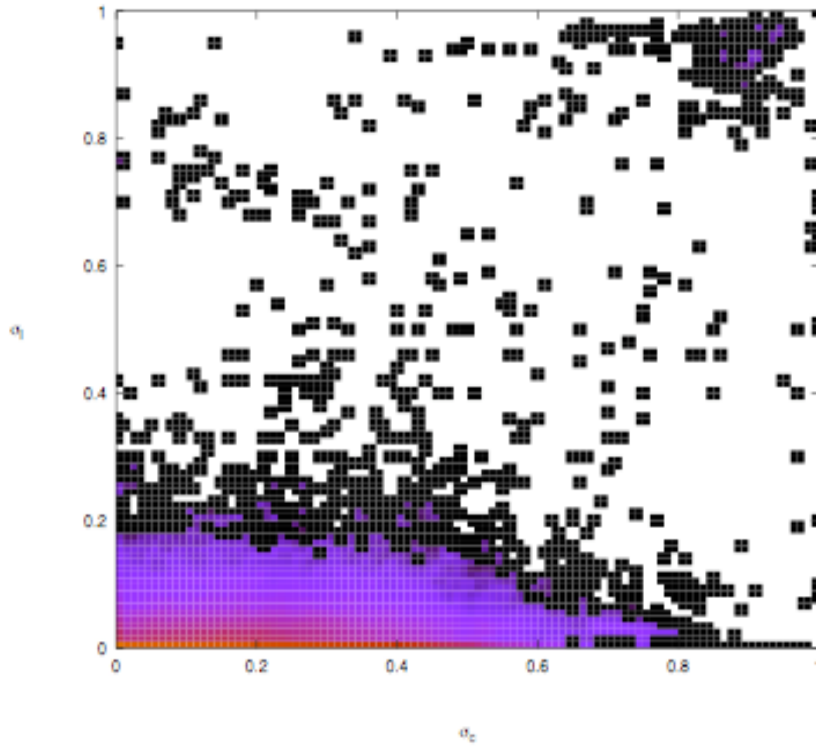
σ_e

Adult

σ_ℓ

log Recall

Precision



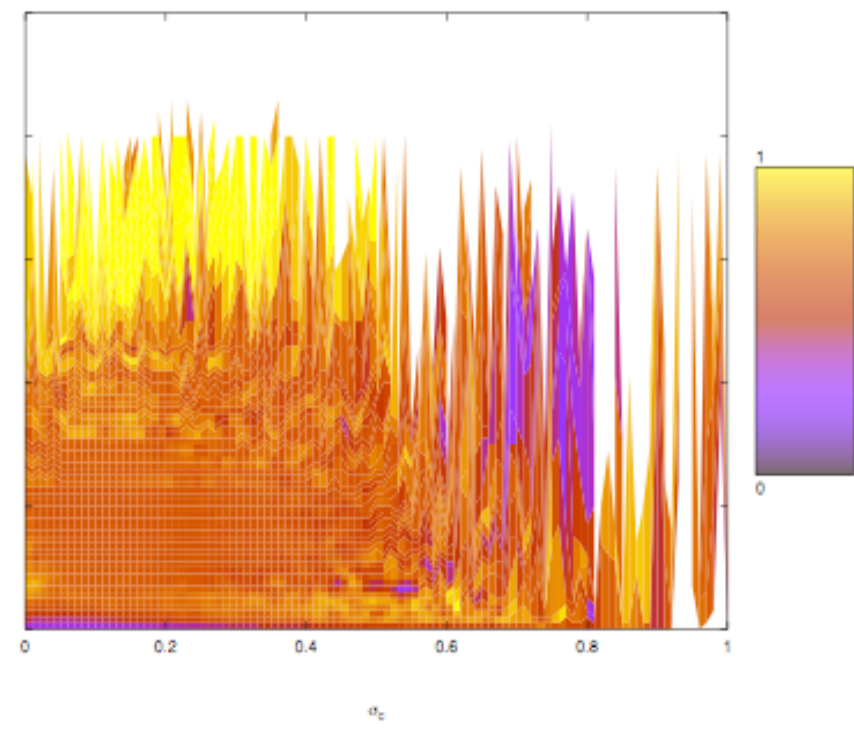
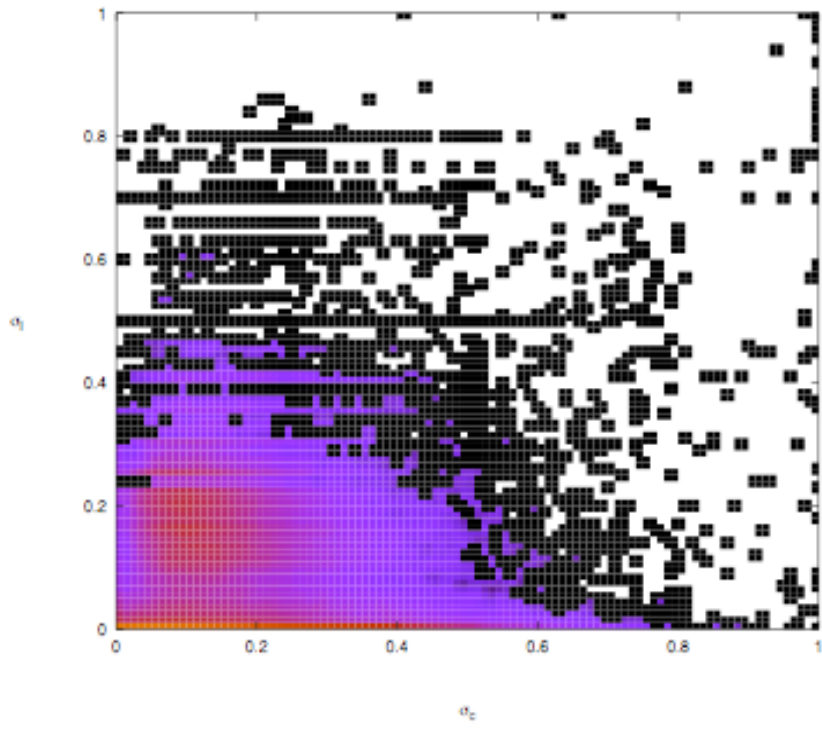
σ_c

Computers

σ_l

log Recall

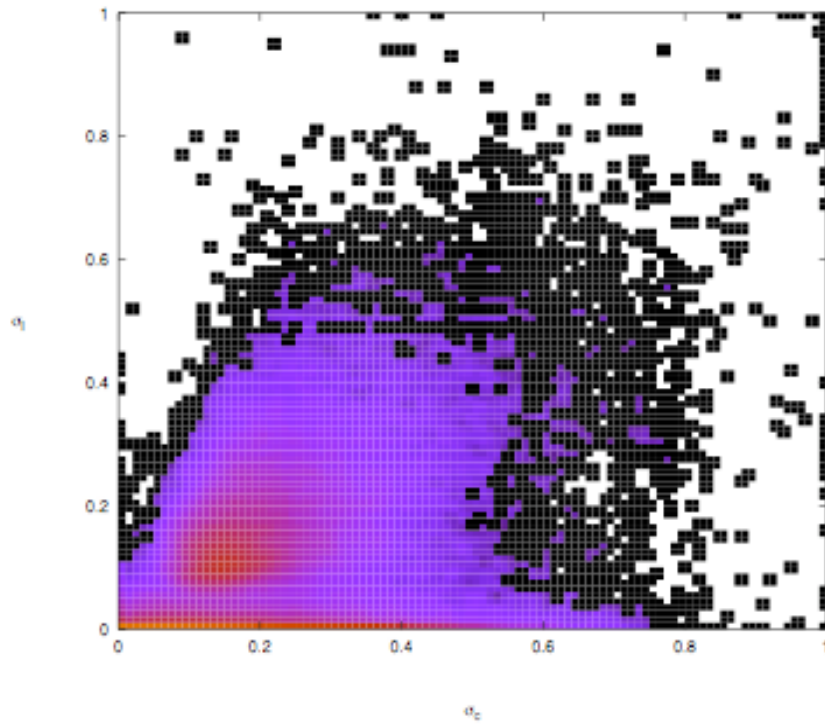
Precision



σ_c

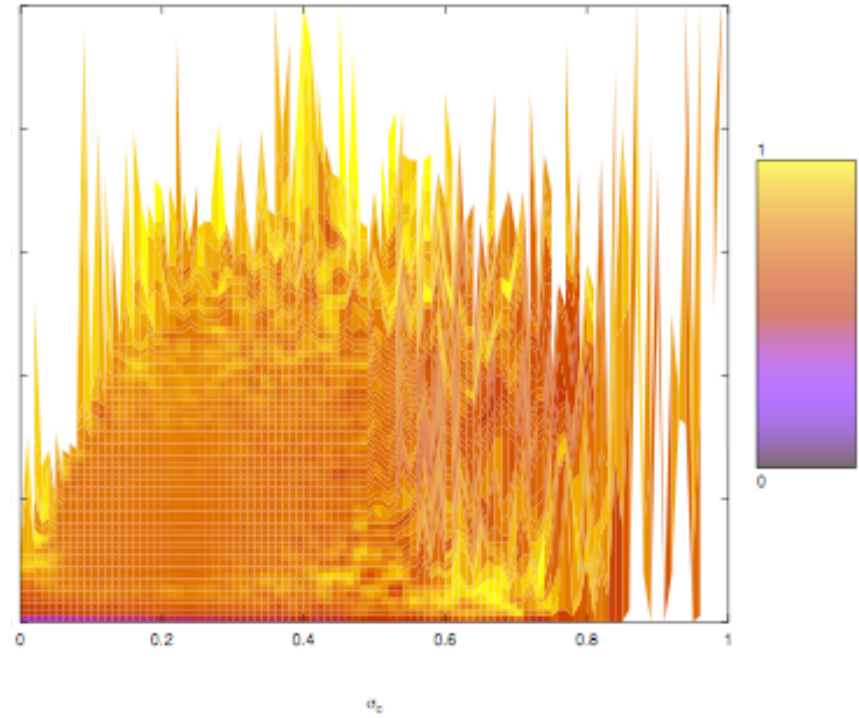
Reference

log Recall



σ_l

Precision



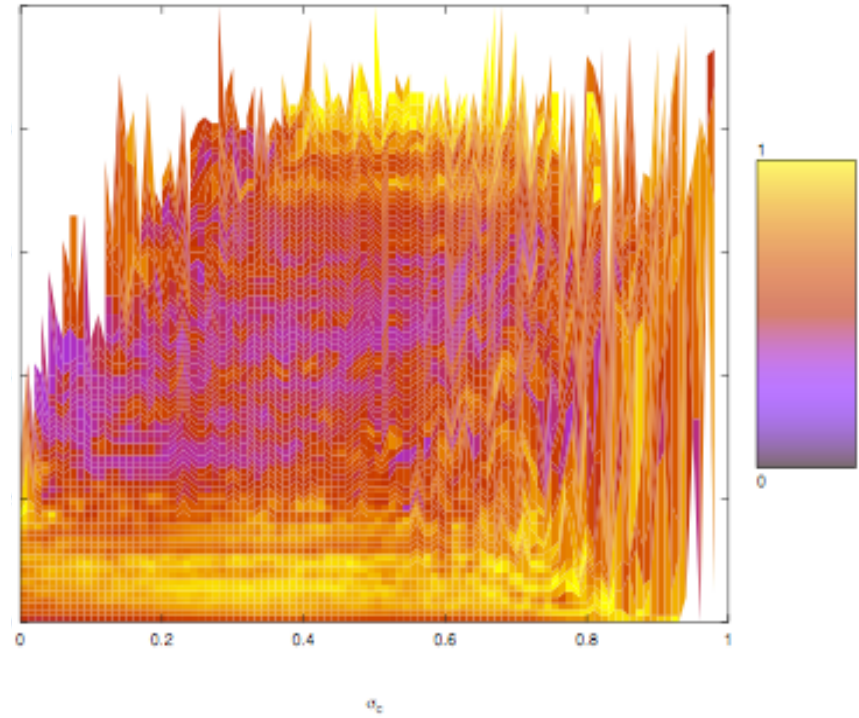
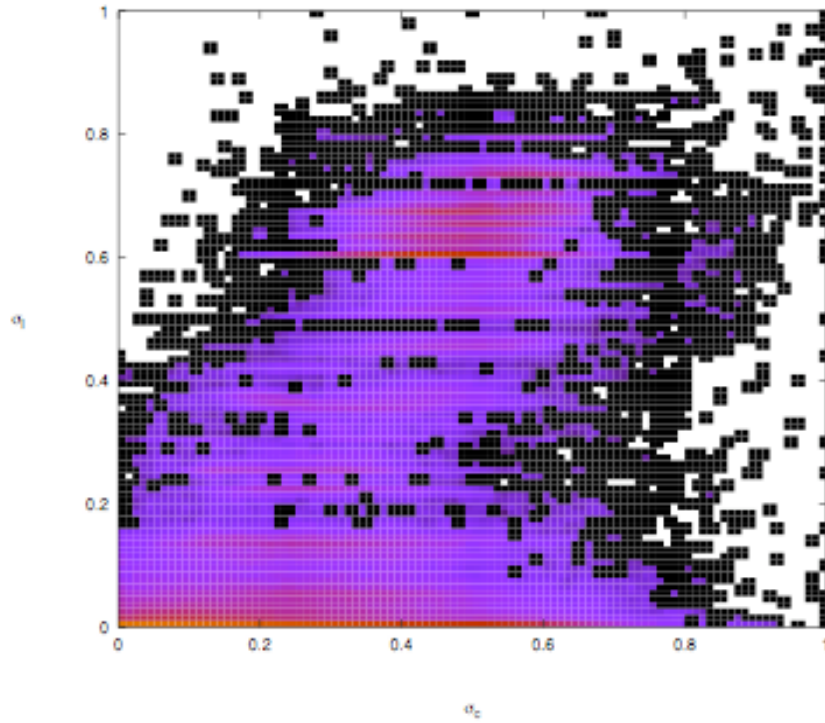
σ_c

Home



log Recall

Precision

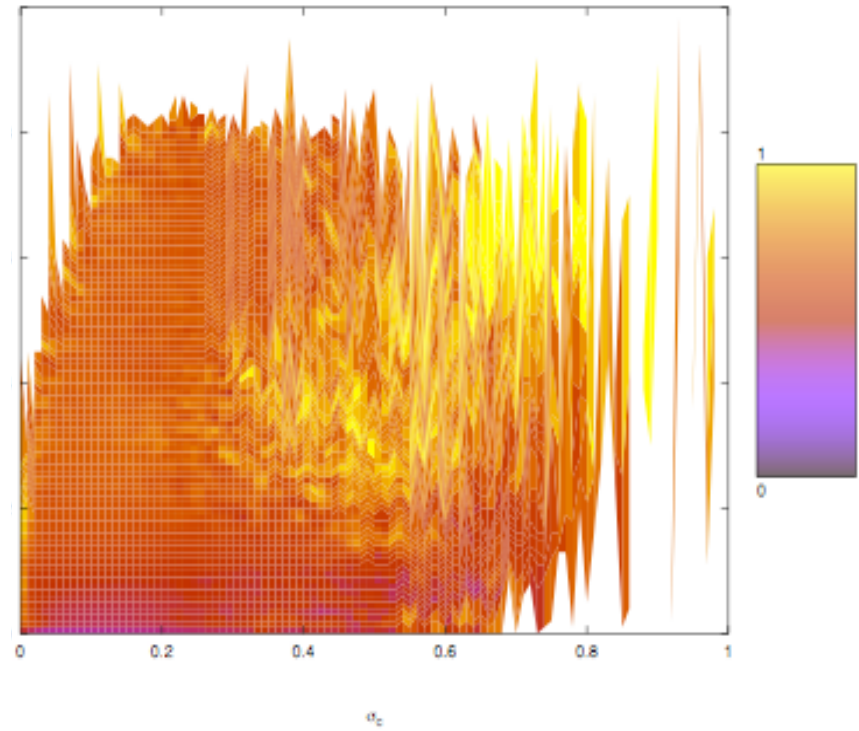
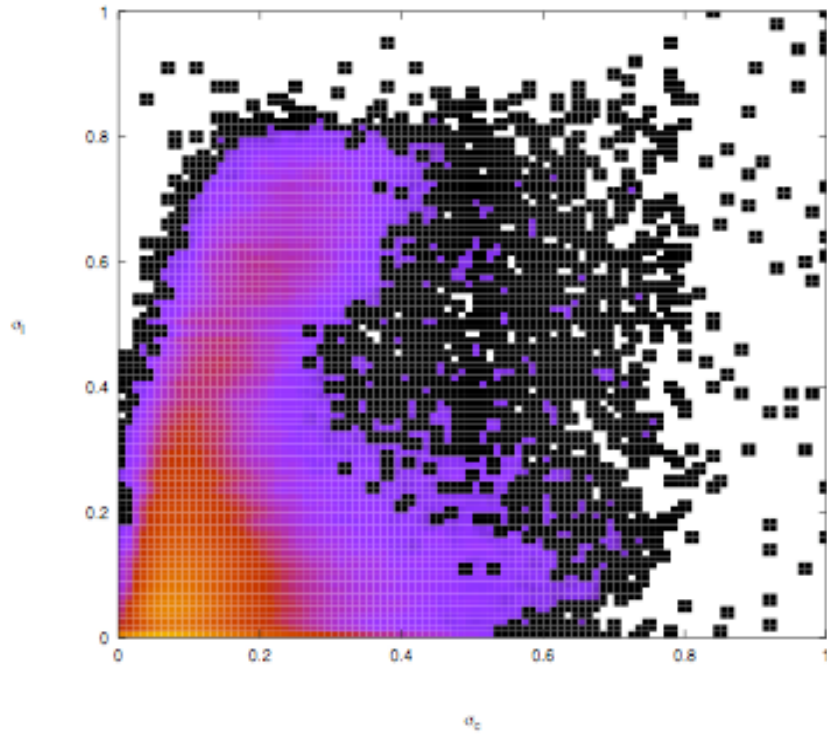


News

σ_ℓ

log Recall

Precision



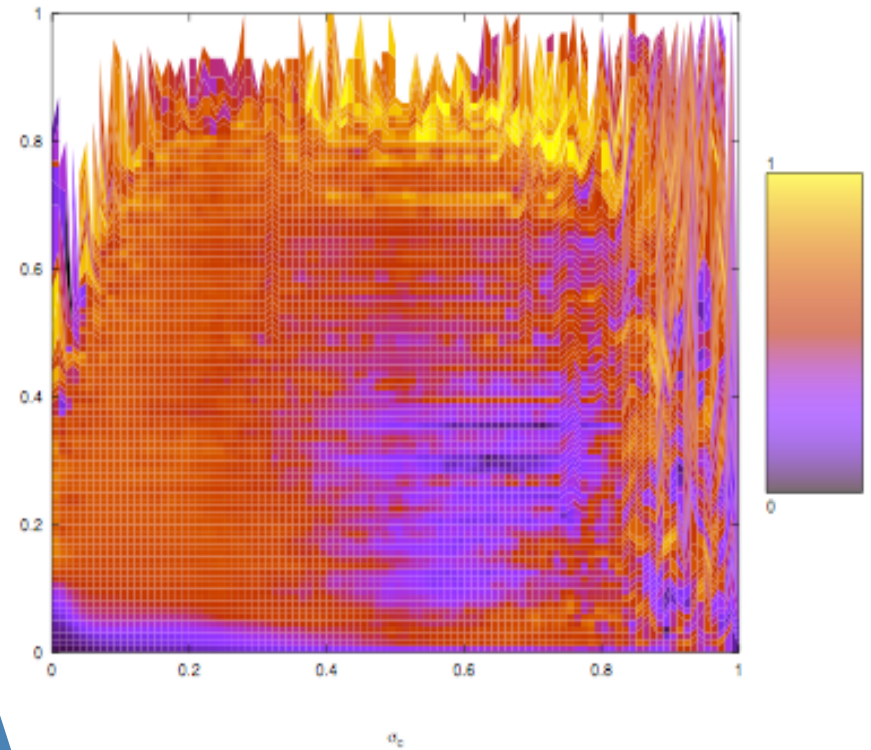
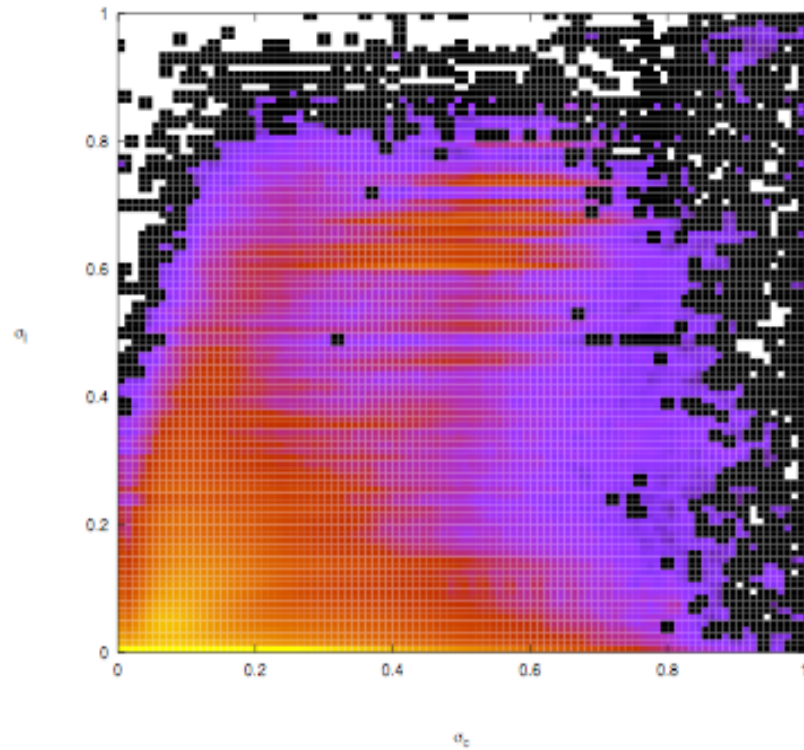
σ_c

All pairs

σ_e

log Recall

Precision

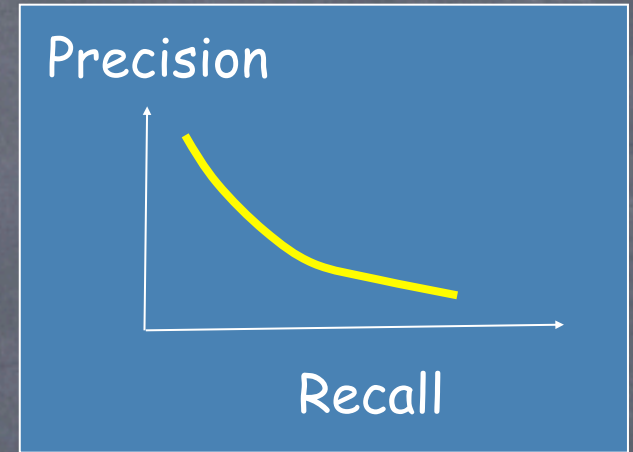


σ_e

Discussion: So what?

So what?

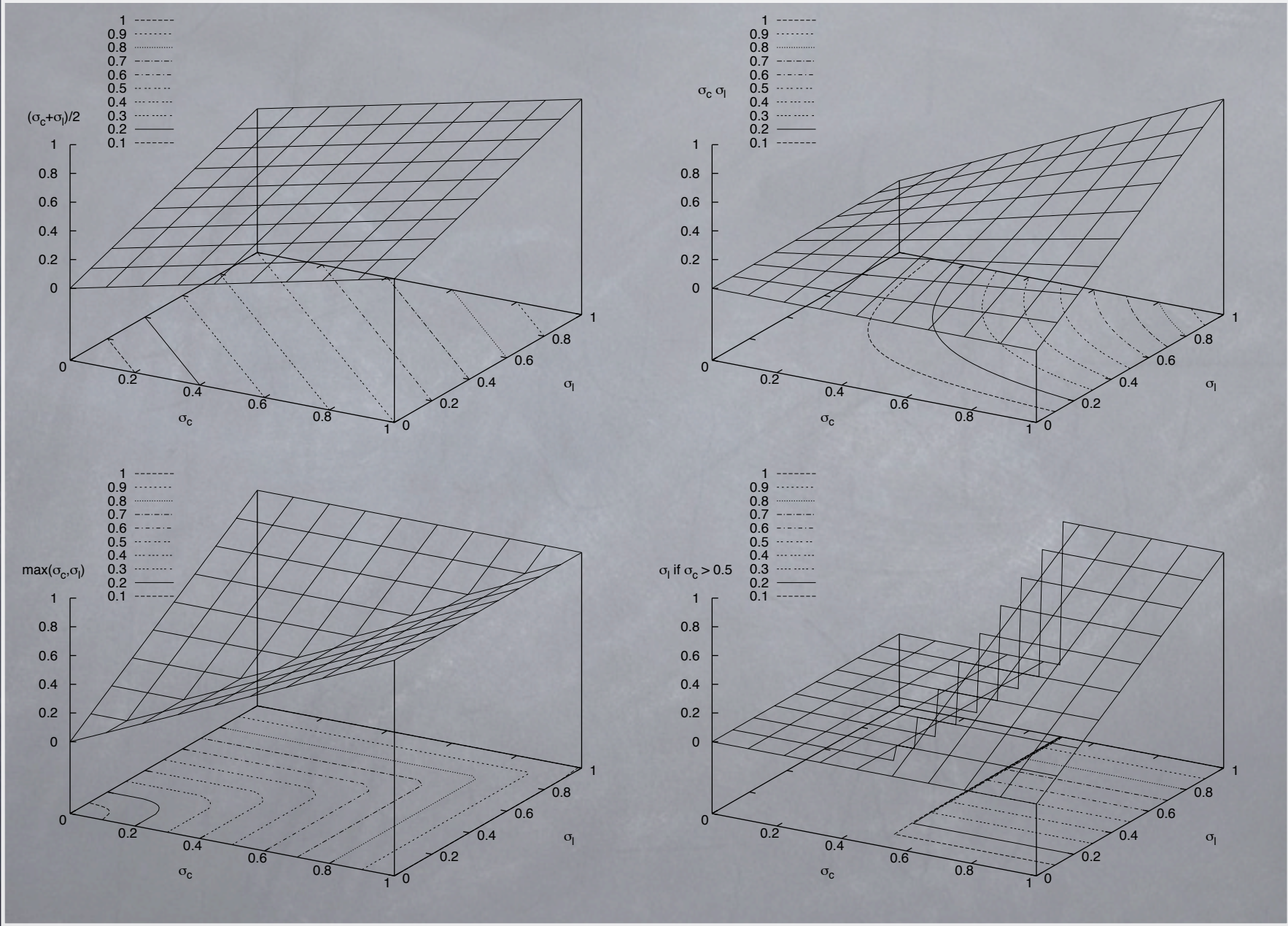
Approximate relevance
by semantic similarity



$$P(f, \beta) = \frac{\sum_{p, q: f(\sigma_c(p, q), \sigma_l(p, q)) \geq \beta} \sigma_s(p, q)}{|\{p, q : f(\sigma_c(p, q), \sigma_l(p, q)) \geq \beta\}|}$$

$$R(f, \beta) = \frac{\sum_{p, q: f(\sigma_c(p, q), \sigma_l(p, q)) \geq \beta} \sigma_s(p, q)}{\sum_{p, q} \sigma_s(p, q)}$$

Rank by combining content and link similarity



All pairs

$$\sigma_l \cdot H(\sigma_c - 0.5)$$

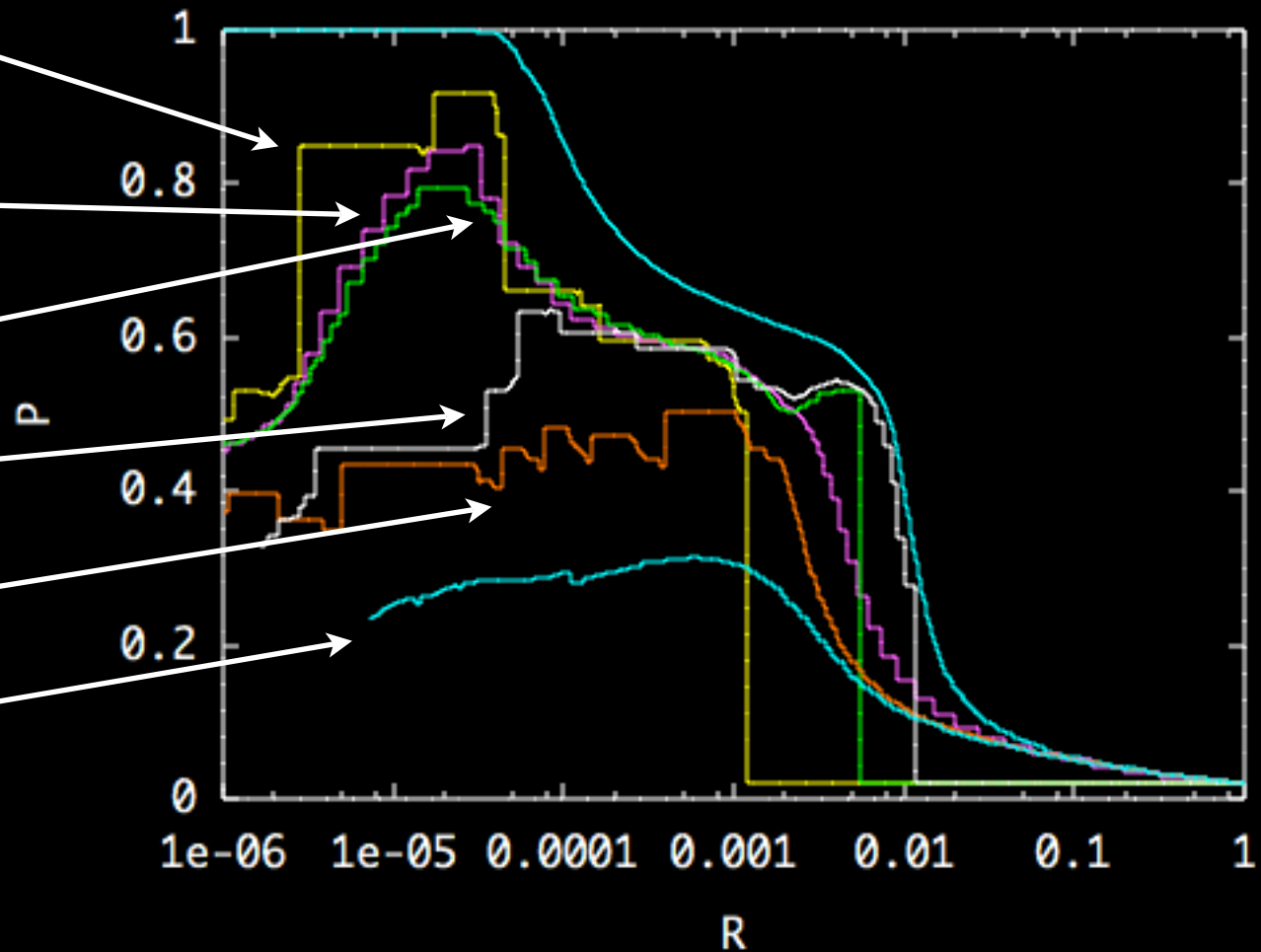
$$\frac{2}{3}\sigma_l + \frac{1}{3}\sigma_c$$

$$\sigma_c \cdot \sigma_l$$

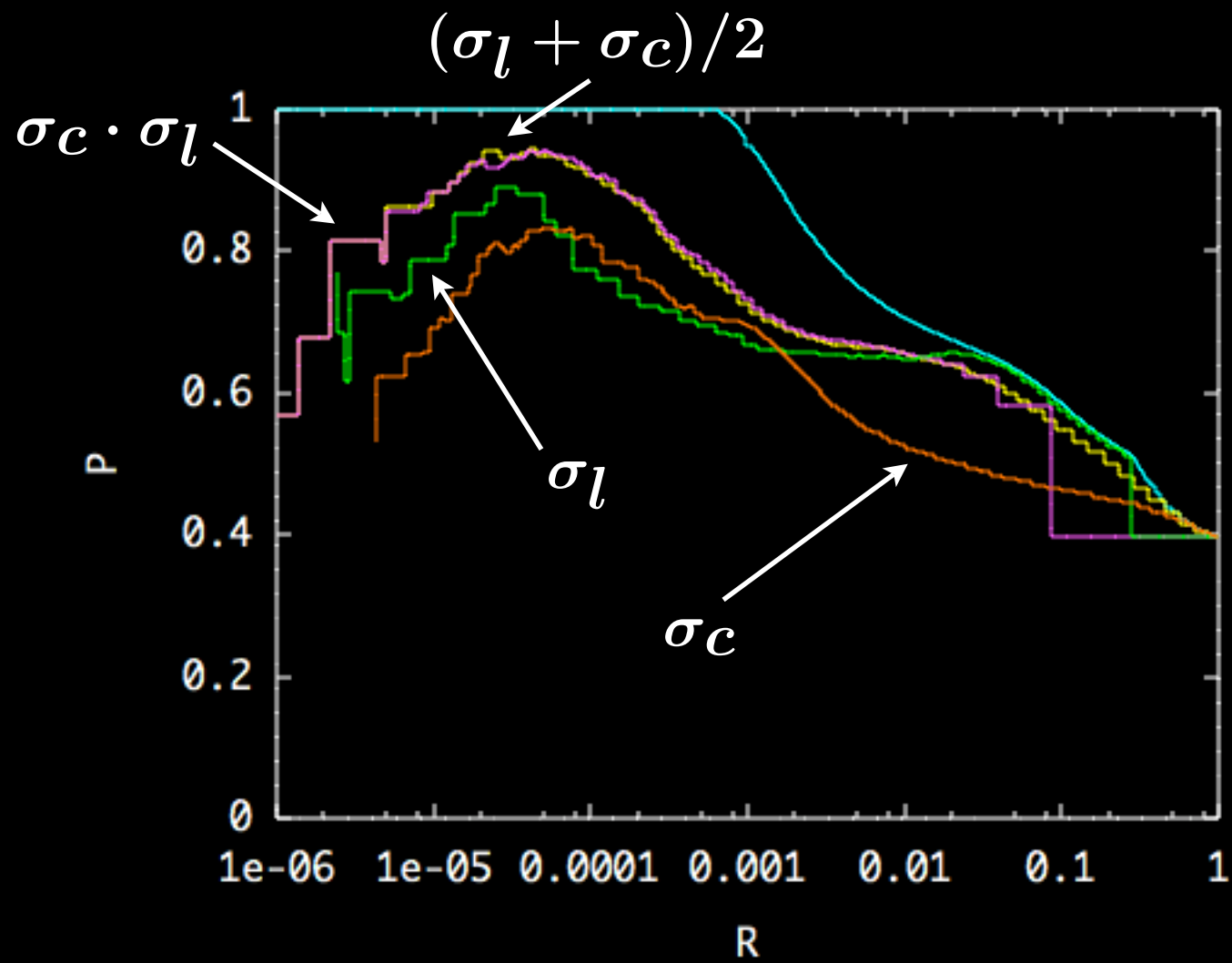
$$\sigma_l$$

$$\max(\sigma_c, \sigma_l)$$

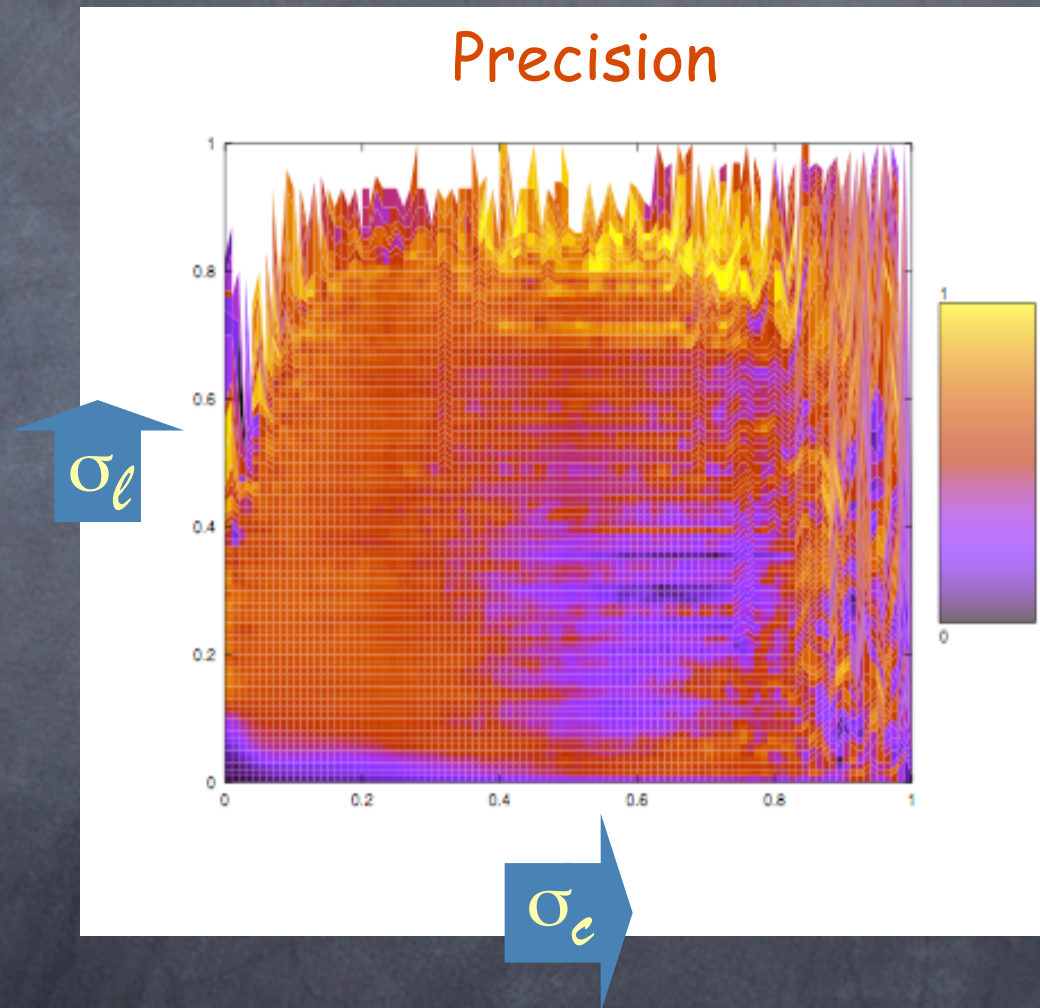
$$\sigma_c$$

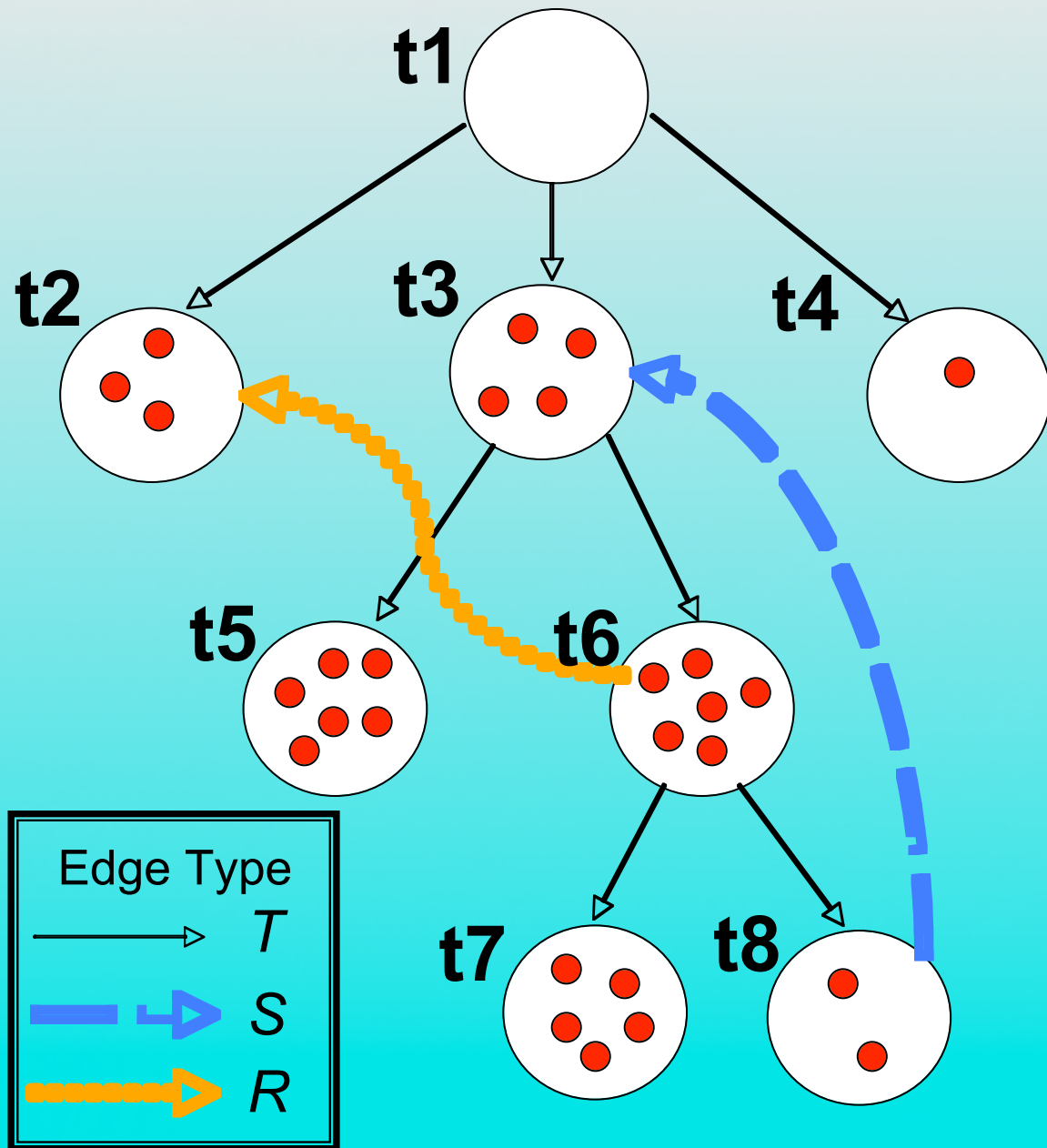


News



What about the "hole"?





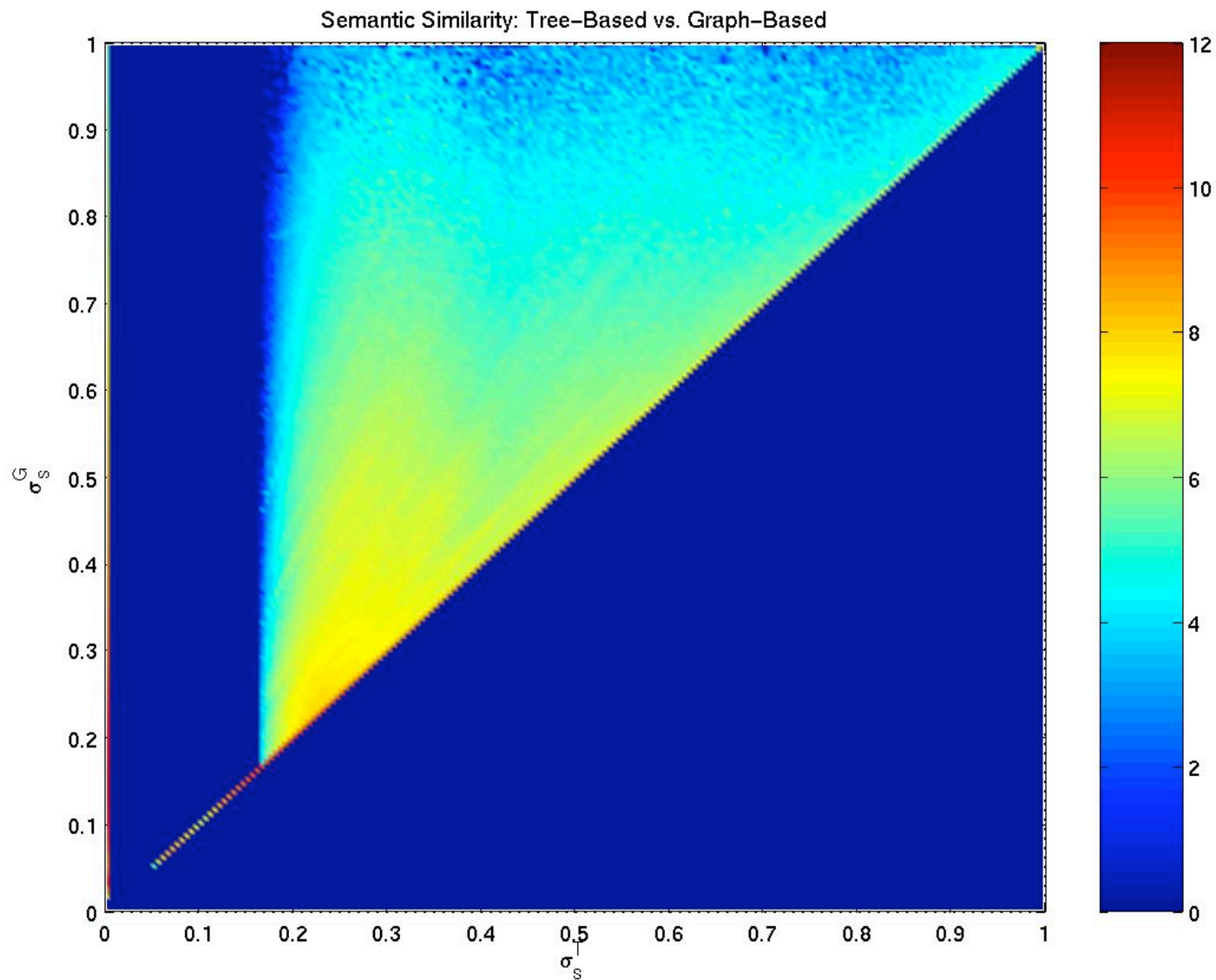
Better semantic similarity measure

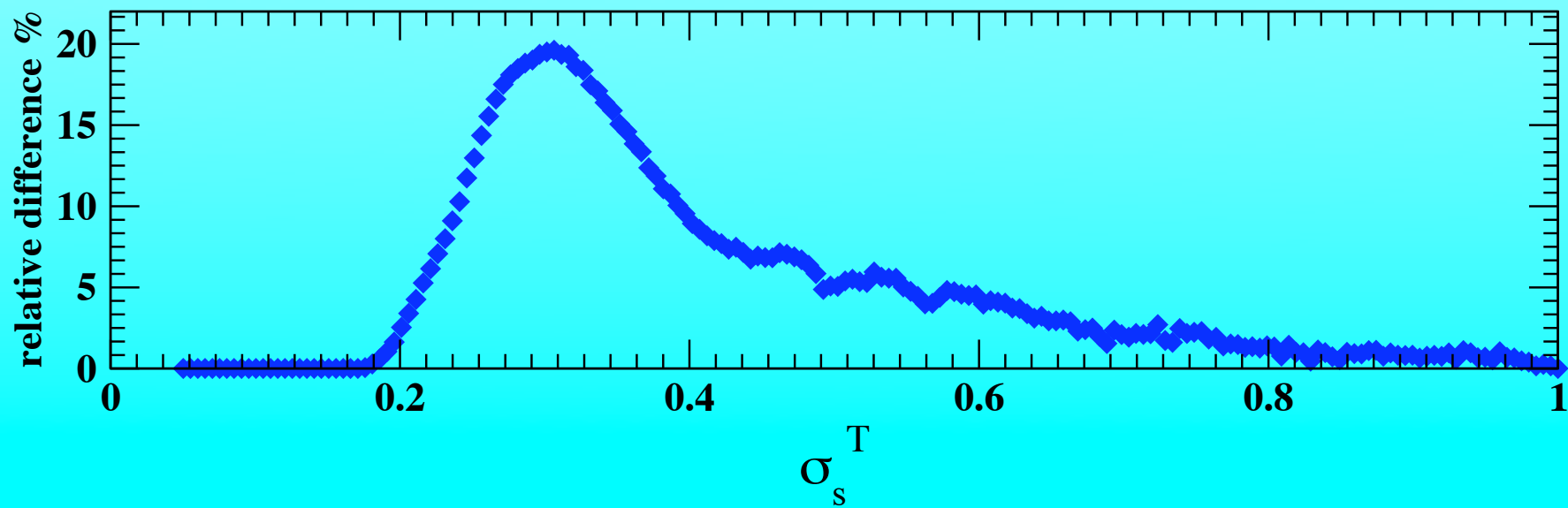
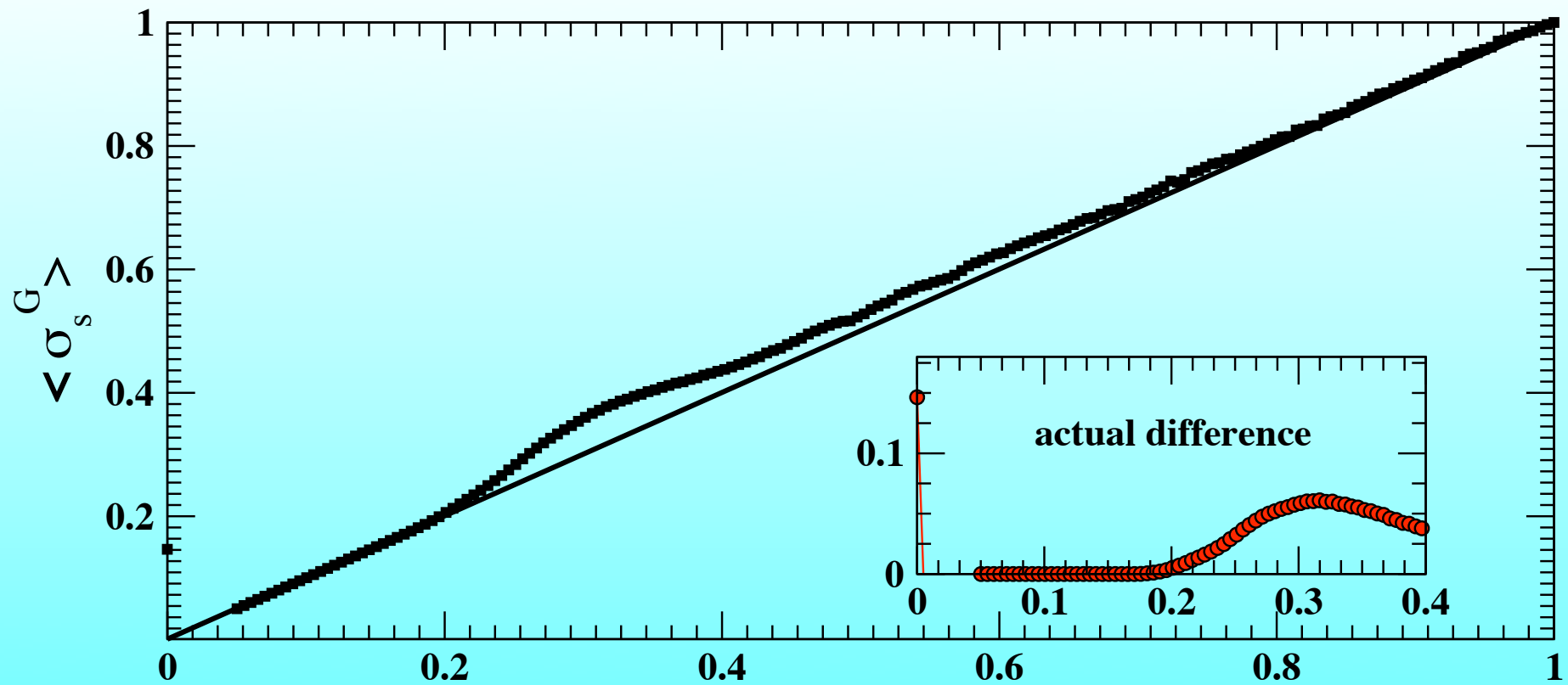
- Work w/ Ana Maguitman & al.
- Include cross-links (symbolic) and see-also links (related)
- Transitive closure of topic graph
- Compute entropy based on fuzzy membership matrix

$$W = T^+ \odot G \odot T^+$$


$$\sigma_s(t_1, t_2) = \max_k \frac{2 \cdot \min(W_{k1}, W_{k2}) \cdot \log \Pr[t_k]}{\log(\Pr[t_1|t_k] \cdot \Pr[t_k]) + \log(\Pr[t_2|t_k] \cdot \Pr[t_k])}$$

Differences





Web Page Relatedness




MuppetsOnline.com


Welcome

[Home](#)
[Shop](#)
[Music](#)
[Pictures](#)
[Sounds](#)

The Eagle Speaks

Hello and welcome to MuppetsOnline.com. This is Sam the Eagle speaking. You can travel around this site merely by clicking the buttons to your left. Home will take you back here. Each section is looked after by a Muppet, so look around as there is lots to do. Visit the Muppet Shop and buy some of our work for yourself or a loved one. It's patriotic and American.




amazon.co.uk
Amazon Recommends:

[The Muppet Christmas Carol](#)
Michael Caine

CLICK HERE TO BUY THIS MOVIE

	mean	stderr
tree	5.7%	0.8%
graph	84.7%	1.8%

"YOUR ONE SOURCE FOR QUALITY FAMILY ENTERTAINMENT"



THE LIVE CAST OF SESAME STREET

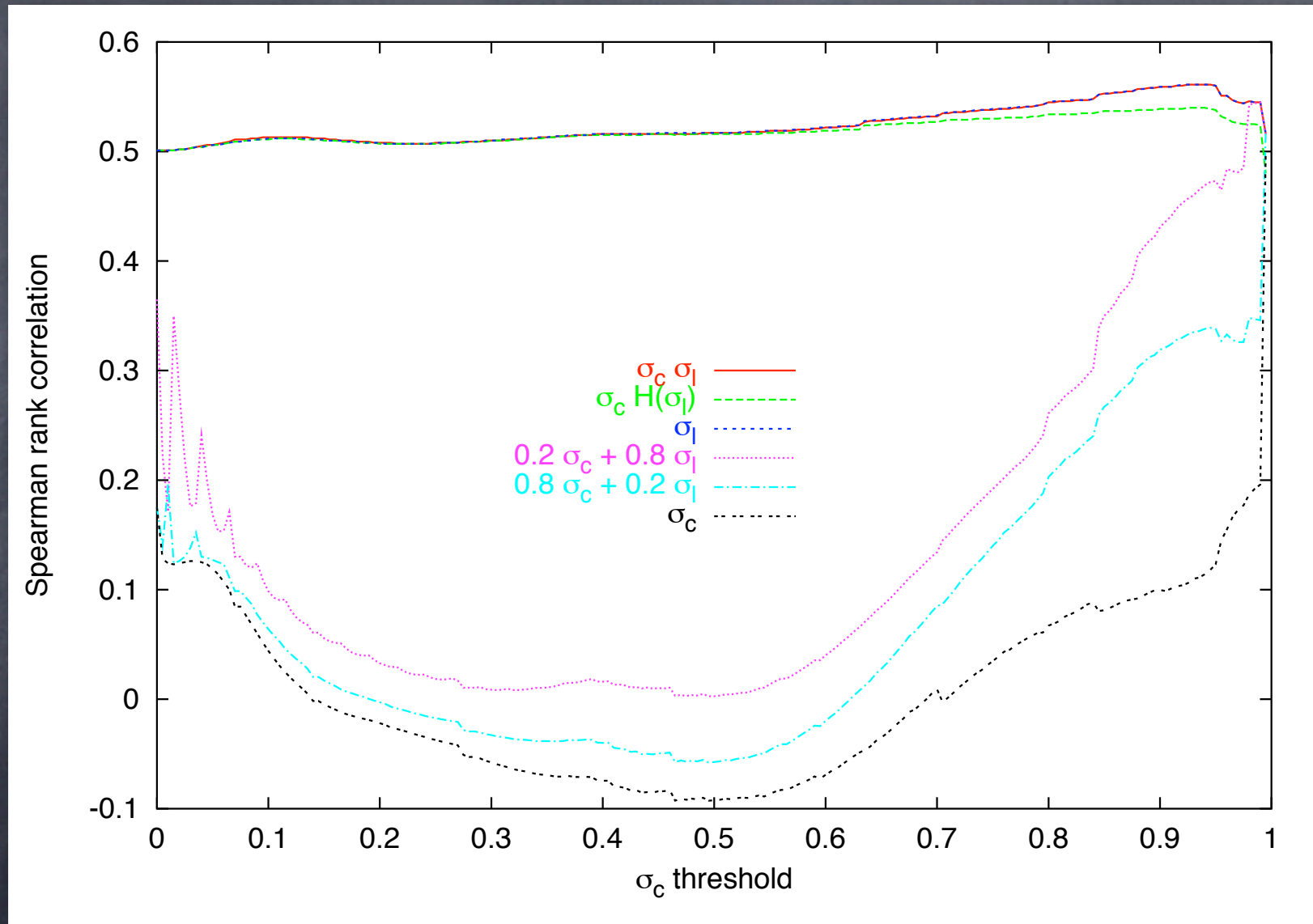
The accolades have all been written many, many times.

Quite simply, there has never been or exists now a television show like SESAME STREET.

SESAME STREET has educated and entertained millions of children all over the world for thirty years! Its name has become synonymous with high quality. Much of the show's success is due to an extremely talented cast of actors.

[ALL ABOUT](#) - [PHOTOS](#) - [UPCOMING EVENTS](#) - [LINKS](#) - [CONTACT US](#)

Combining content & links



Discussion: Is content
really so bad? Why?

Outline

- ✓ Mapping
 - > Topical locality
 - > Content, link, and semantic topologies in the Web
- ◉ Modeling
 - > How the Web evolves and why content matters
 - > Consequences for navigation and crawling
- ◉ Mining
 - > Topical Web crawlers
 - > Adaptive, intelligent crawling techniques
- ◉ Mingling
 - > Social Web search & recommendation
 - > Distributed collaborative peer search

Preferential attachment

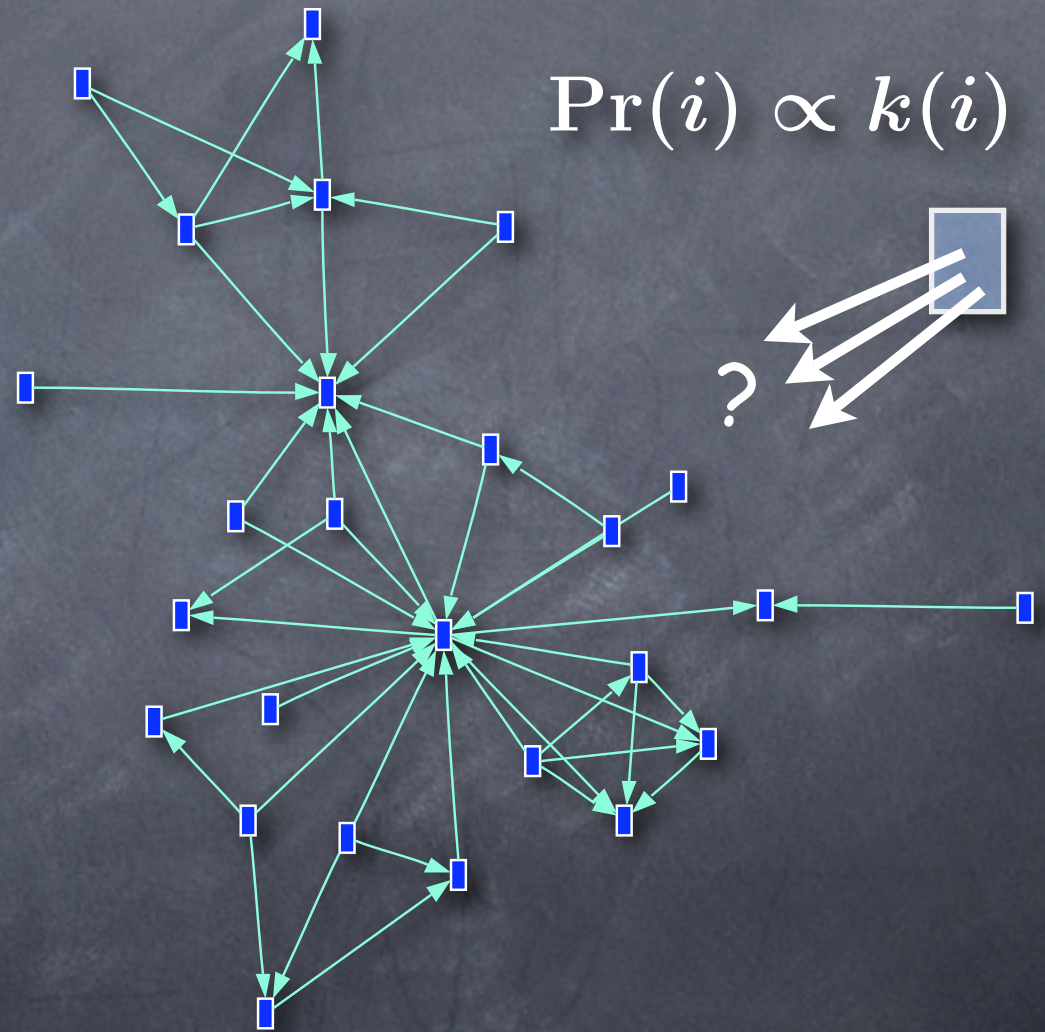
“BA” model

(Barabasi & Albert 1999,
de Solla Price 1976)

- > At each step t
add new page p
- > Create m new
links from p to i
($i < t$)

Rich-get-richer

$$\Pr(i) = \frac{k(i)}{mt} \implies \Pr(k) \sim k^{-\gamma}$$



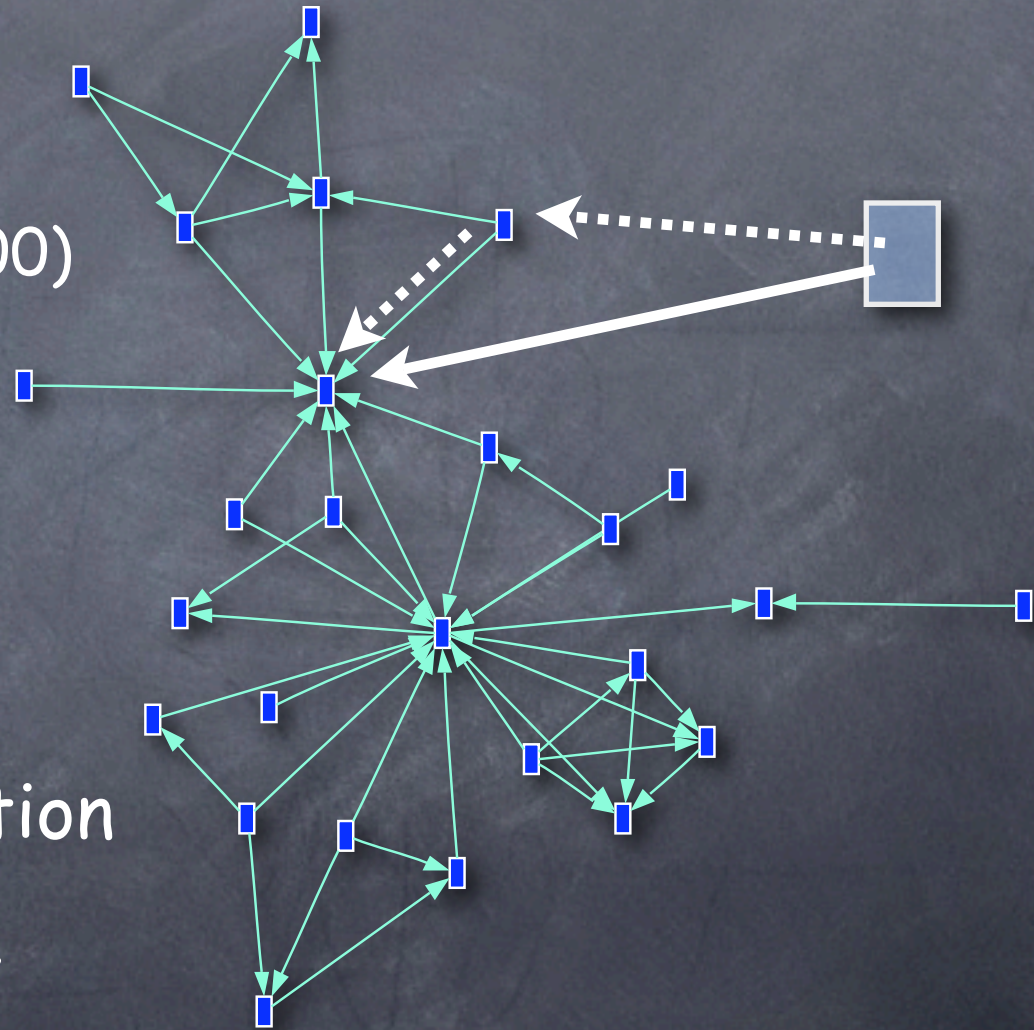
Other growth models

Web copying

(Kleinberg, Kumar & al 1999, 2000)

$$\Pr(i) \propto \Pr(j) \cdot \Pr(j \rightarrow i)$$

- same indegree distribution
- no need to know degree



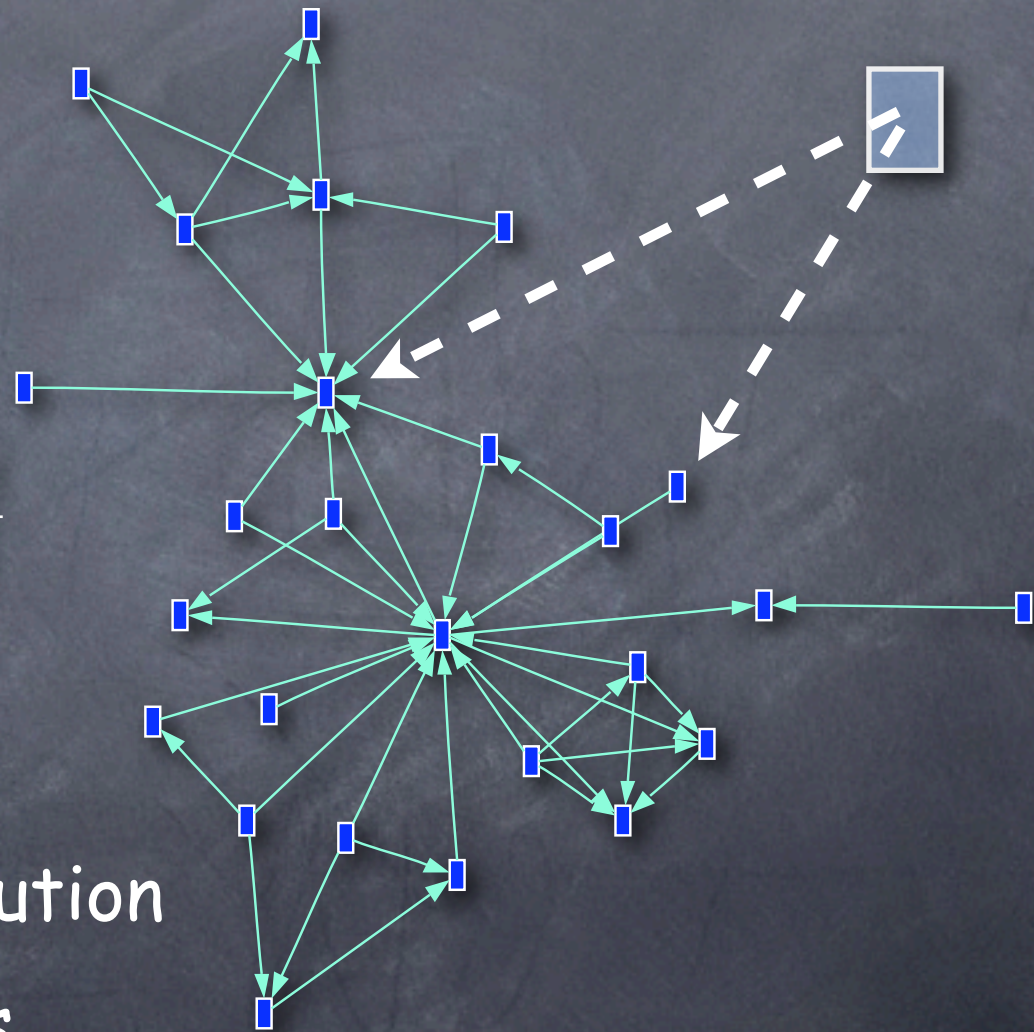
Other growth models

Random mixture

(Pennock & al. 2002,
Cooper & Frieze 2001,
Dorogovtsev & al 2000)

$$\Pr(i) \propto \psi \cdot \frac{1}{t} + (1 - \psi) \cdot \frac{k(i)}{mt}$$

- winners don't take all
- general indegree distribution
- fits non-power-law cases



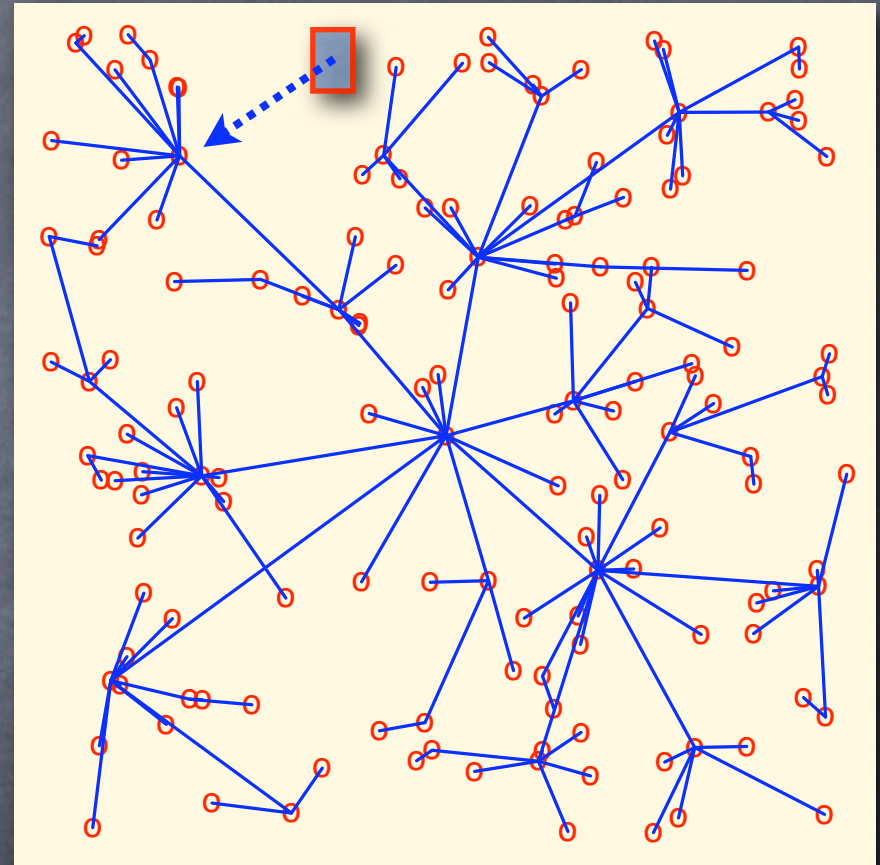
Other growth models

Mixture with Euclidean distance

(HOT: Fabrikant, Koutsoupias
& Papadimitriou 2002)

$$i = \arg \min(\phi r_{it} + g_i)$$

- tradeoff between centrality and geometric locality
- fits power-law in certain critical trade-off regimes

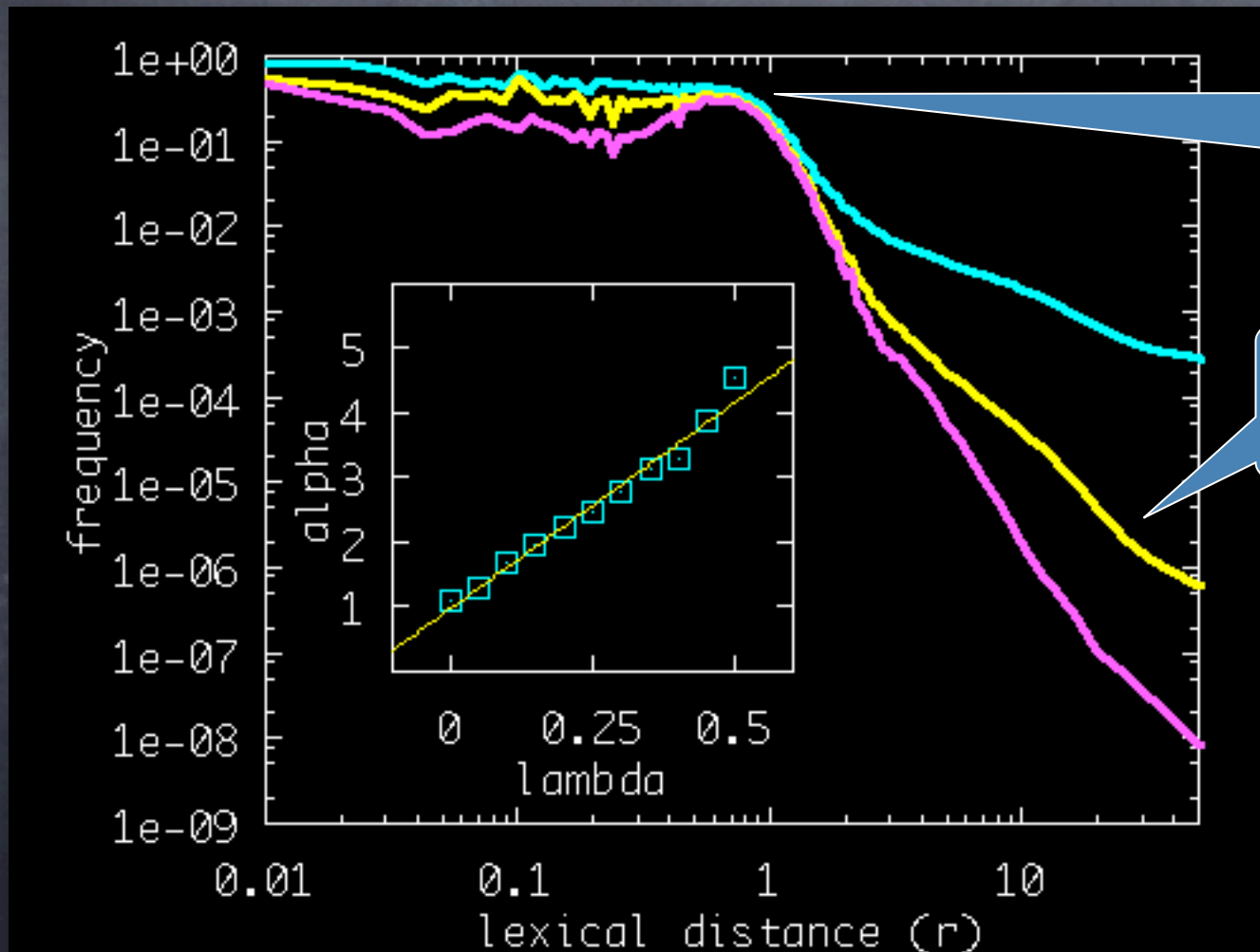


Discussion:
What about content?

Link probability vs lexical distance

$$r = 1/\sigma_c - 1$$

$$\Pr(\lambda | \rho) = \frac{|(p,q) : r = \rho \wedge \sigma_l > \lambda|}{|(p,q) : r = \rho|}$$



Phase
transition
 ρ^*

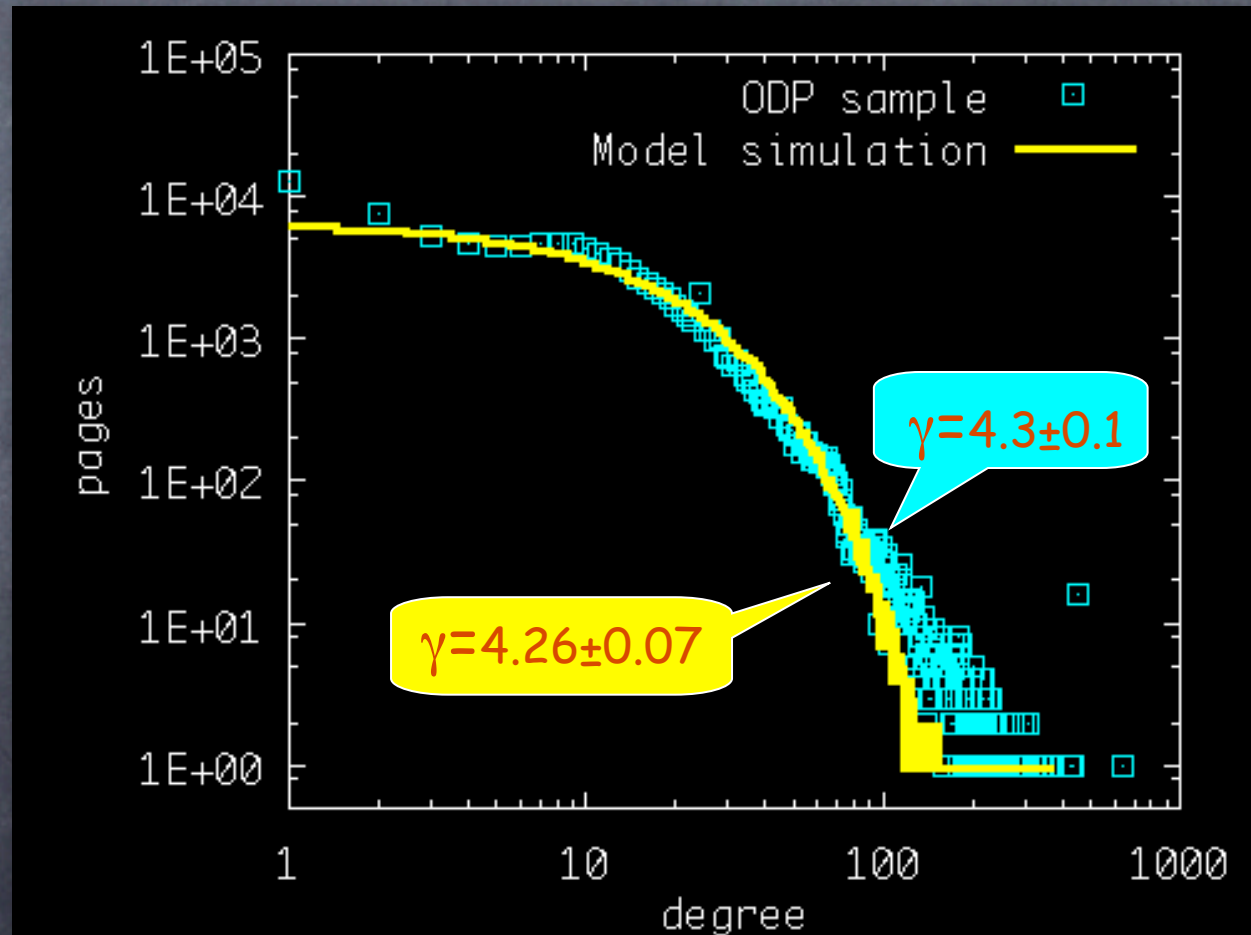
Power law tail
 $\Pr(\lambda | \rho) \sim \rho^{-\alpha(\lambda)}$

*Proc. Natl. Acad.
Sci. USA 99(22):
14014-14019, 2002*

Local content-based growth model

$$\Pr(p_t \rightarrow p_{i < t}) = \begin{cases} \frac{k(i)}{mt} & \text{if } r(p_i, p_t) < \rho^* \\ c[r(p_i, p_t)]^{-\alpha} & \text{otherwise} \end{cases}$$

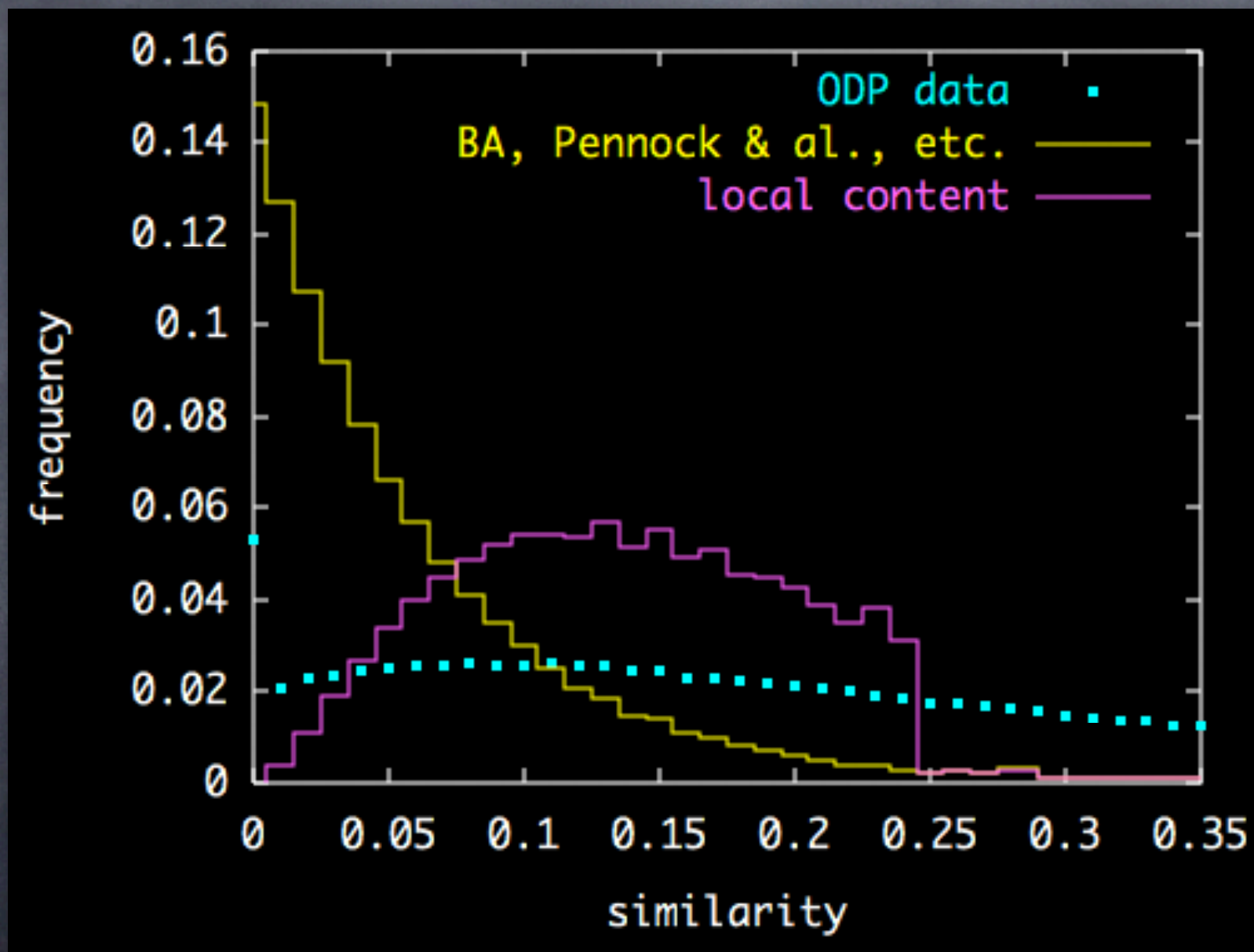
- Similar to preferential attachment (BA)
- Use degree info (popularity/ importance) only for nearby (similar/ related) pages



So, many models can predict degree distributions...

- 👁 Which is "right" ?
- 👁 Need an independent observation (other than degree) to validate models
- 👁 Distribution of content similarity across linked pairs

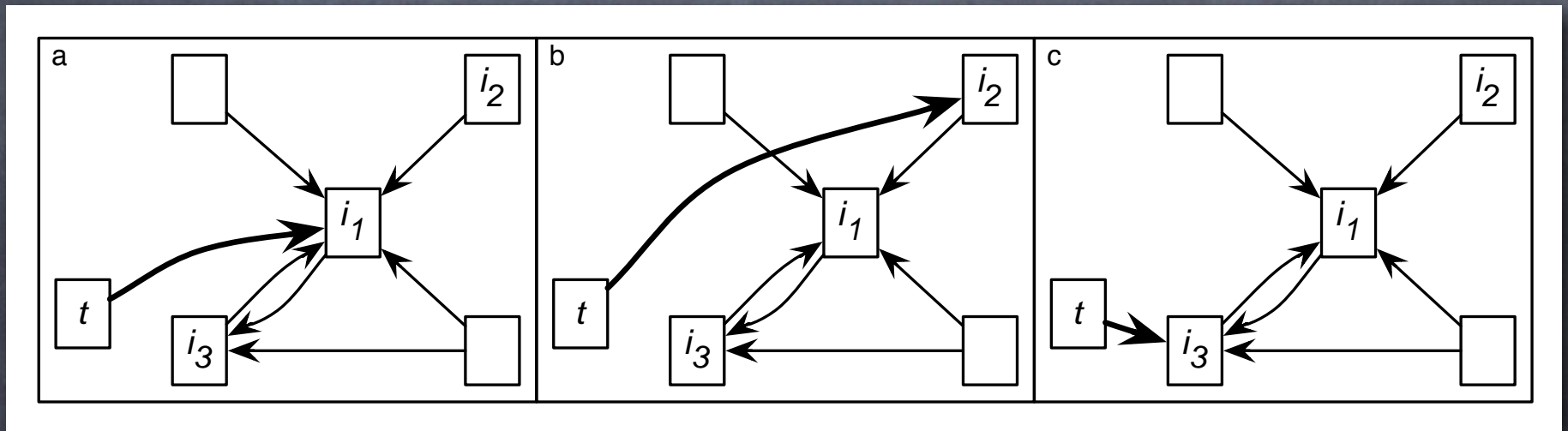
None of these models is right!



Back to the mixture model

$$\Pr(i) \propto \psi \cdot \frac{1}{t} + (1 - \psi) \cdot \frac{k(i)}{mt}$$

degree-uniform mixture



Bias choice by content similarity instead
of uniform distribution

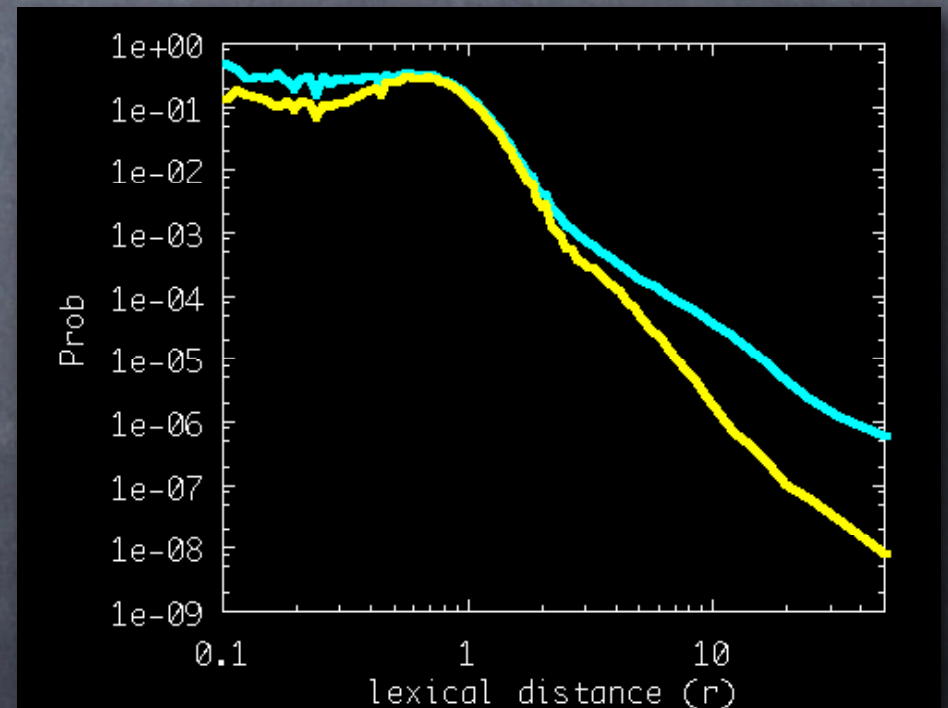
Degree-similarity mixture model

$$\text{Pr}(i) \propto \psi \cdot \hat{\text{Pr}}(i) + (1 - \psi) \cdot \frac{k(i)}{mt}$$

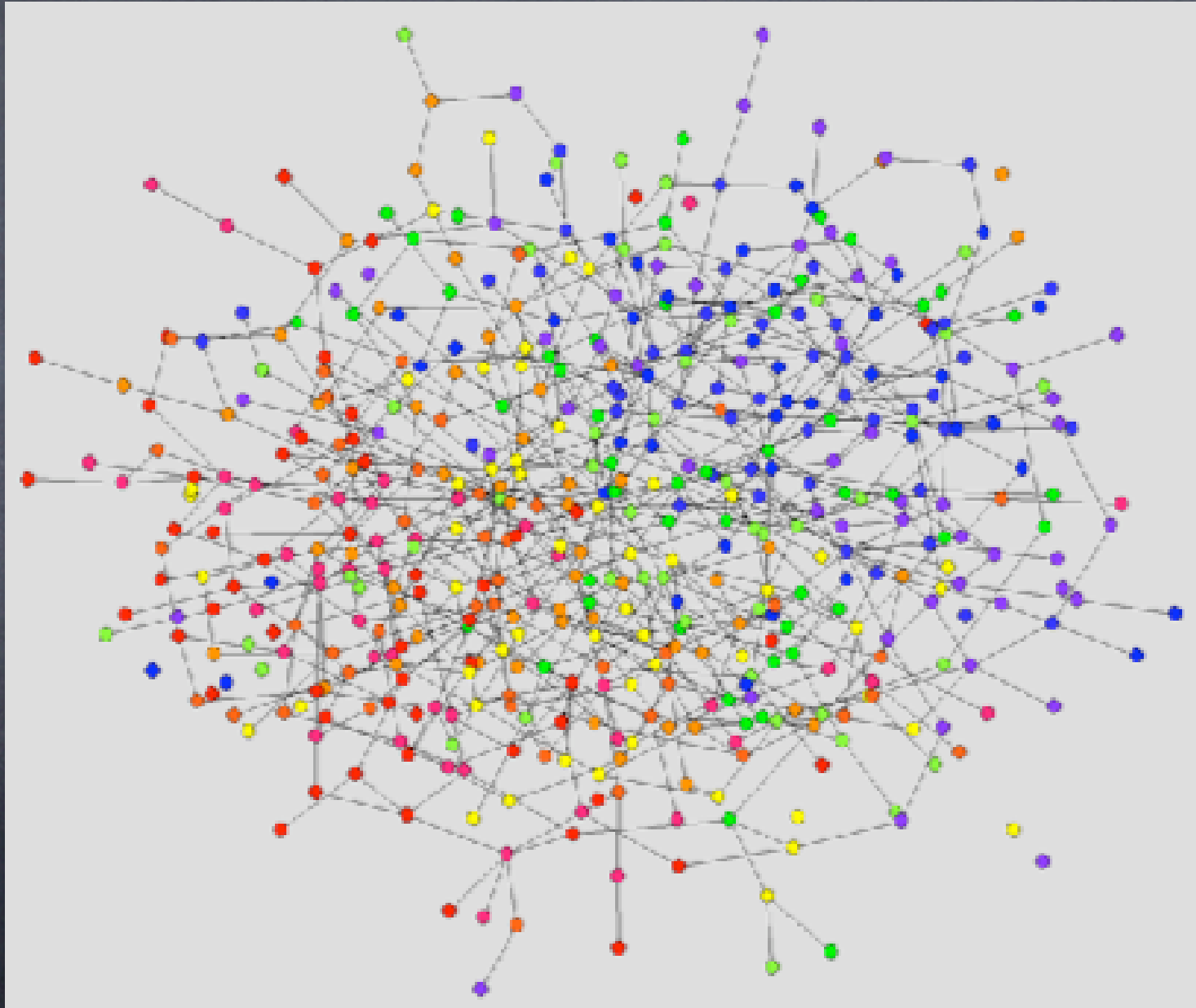


$$\hat{\text{Pr}}(i) \propto [r(i, t)]^{-\alpha}$$

$$\psi = 0.2, \alpha = 1.7$$

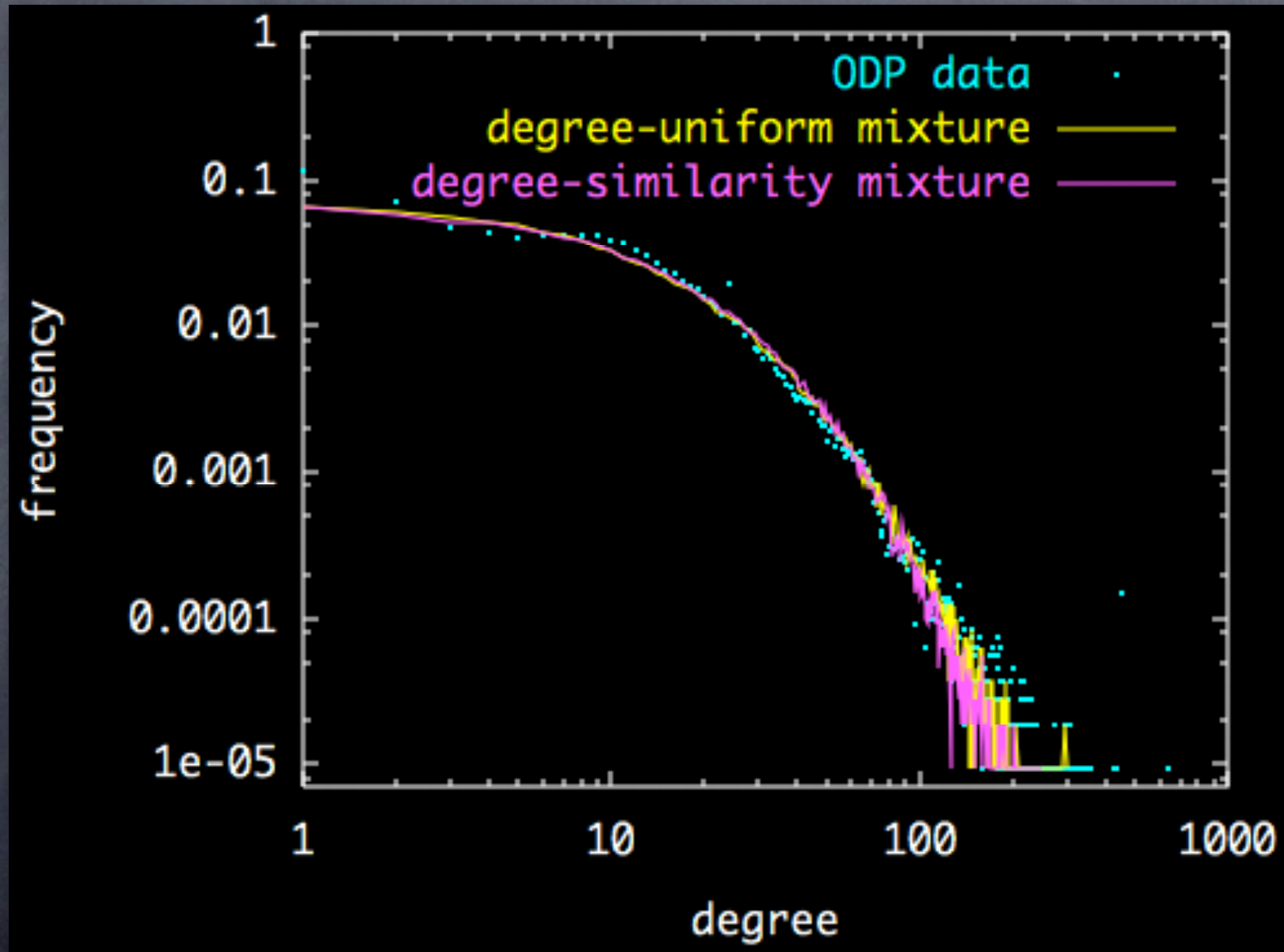


Build it...

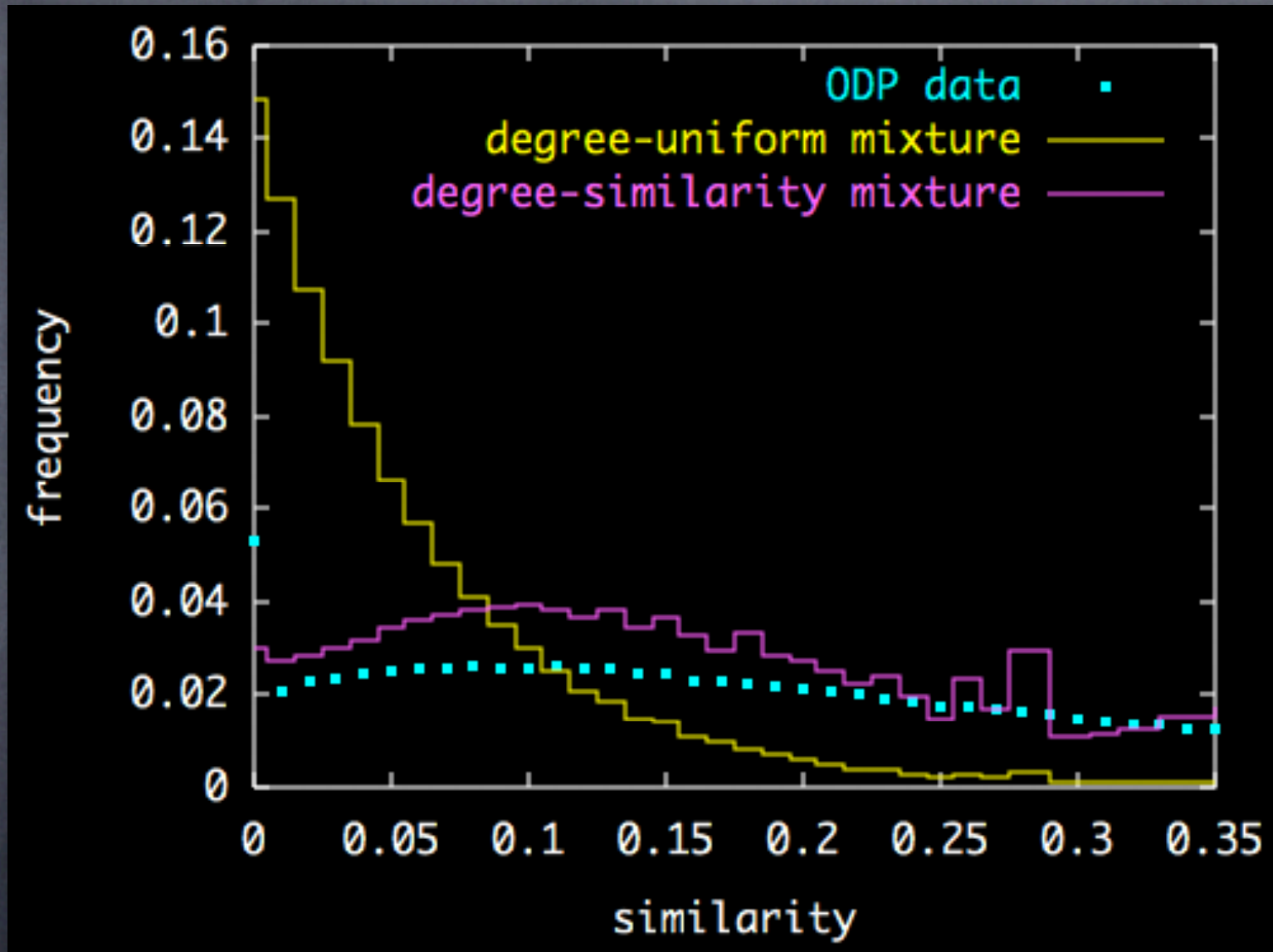


(M.M.)

Both mixture models get the degree distribution right...

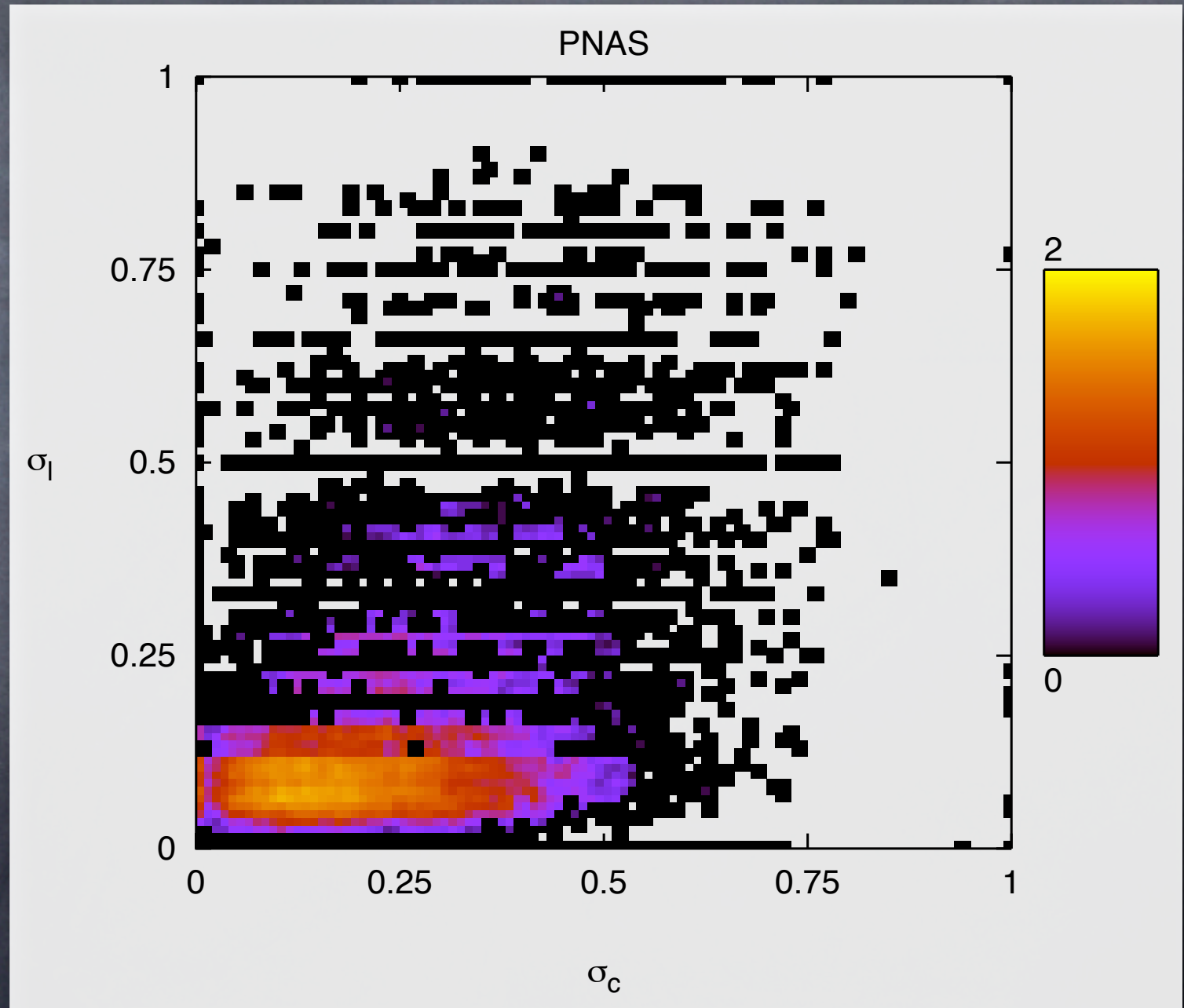


...but the degree-similarity mixture model predicts the similarity distribution better

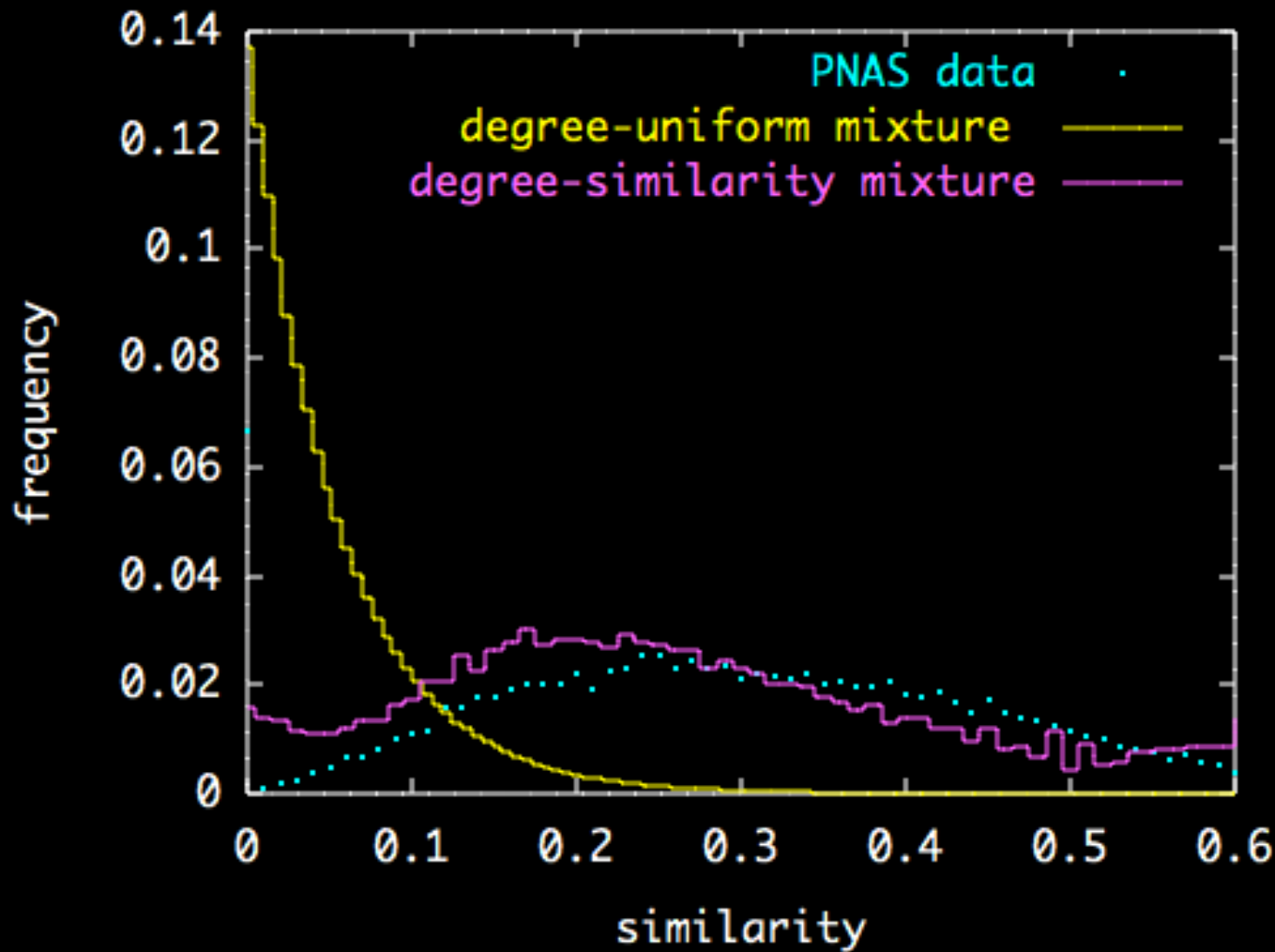
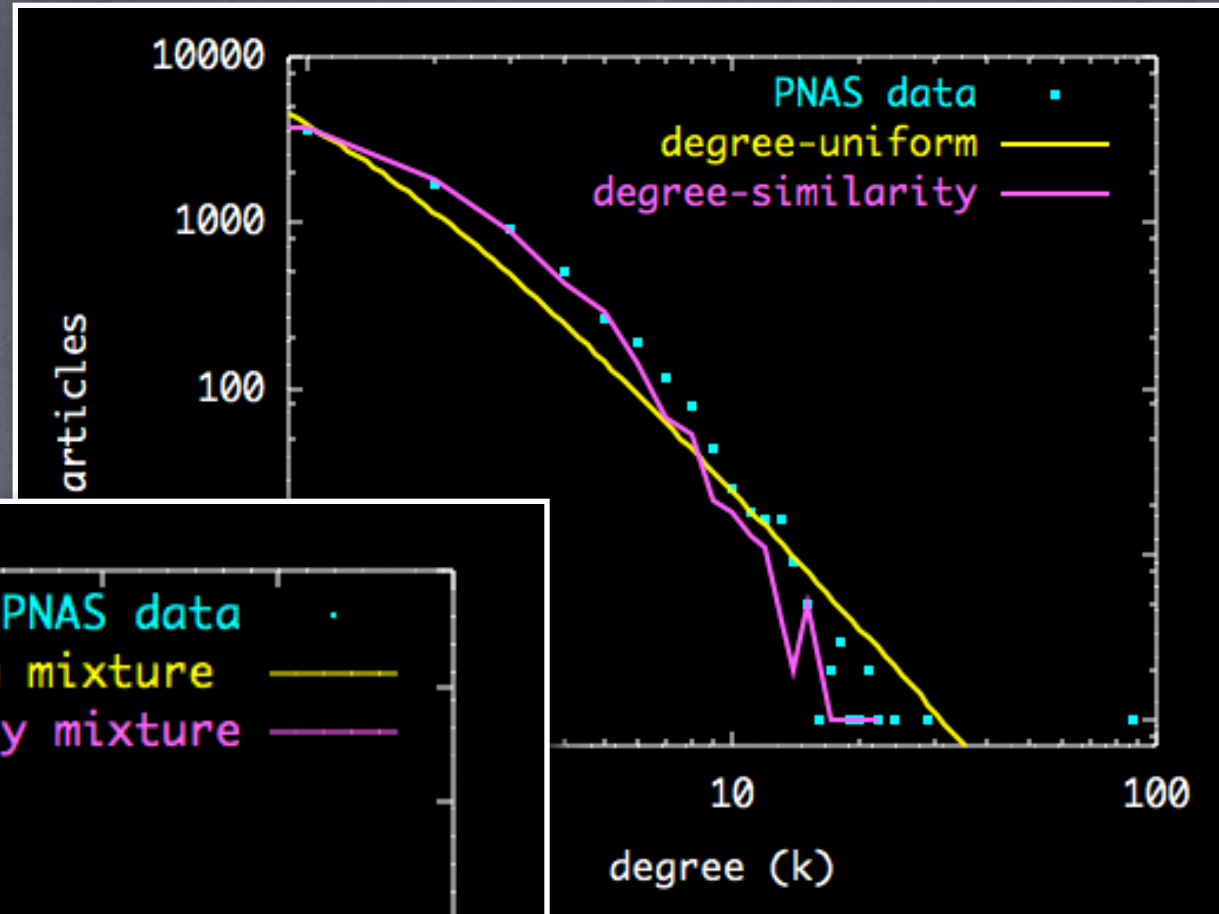


Citation networks

15,785
articles
published
in PNAS
between
1997 and
2002



Citation networks



Discussion

- What questions remain open?
- What other factors should we consider?

Open Questions

- ① Understand distribution of content similarity across all pairs of pages
 - ① Using TF: exponential $\Pr(\sigma) \sim \alpha^{-\beta\sigma}$
 - ① Using TF-IDF: power law $\Pr(\sigma) \sim \alpha\sigma^{-\beta}$
- ① Growth model to explain co-evolution of both link topology and content similarity
- ① The role of search engines (keynote)

Efficient crawling algorithms?

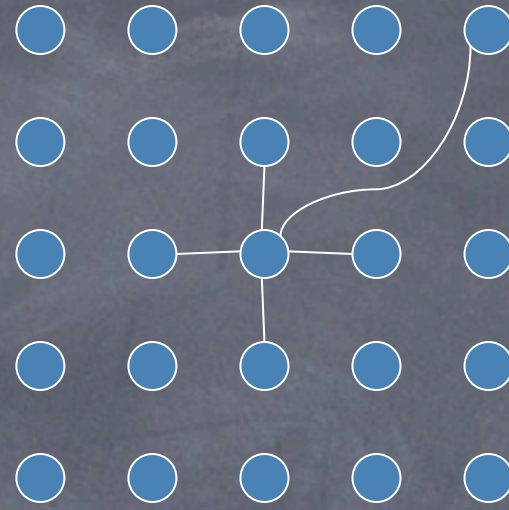
Theory: since the Web is a small world network, or has a scale free degree distribution, **short paths exist** between any two pages:

- ~ $\log N$ (Barabasi & Albert 1999)
- ~ $\log N / \log \log N$ (Bollobas 2001)

Practice: **can't find them!**

- Greedy algorithms based on location in geographical small world networks: ~ $\text{poly}(N)$ (Kleinberg 2000)
- Greedy algorithms based on degree in power law networks: ~ N (Adamic, Huberman & al. 2001)

Exception # 1



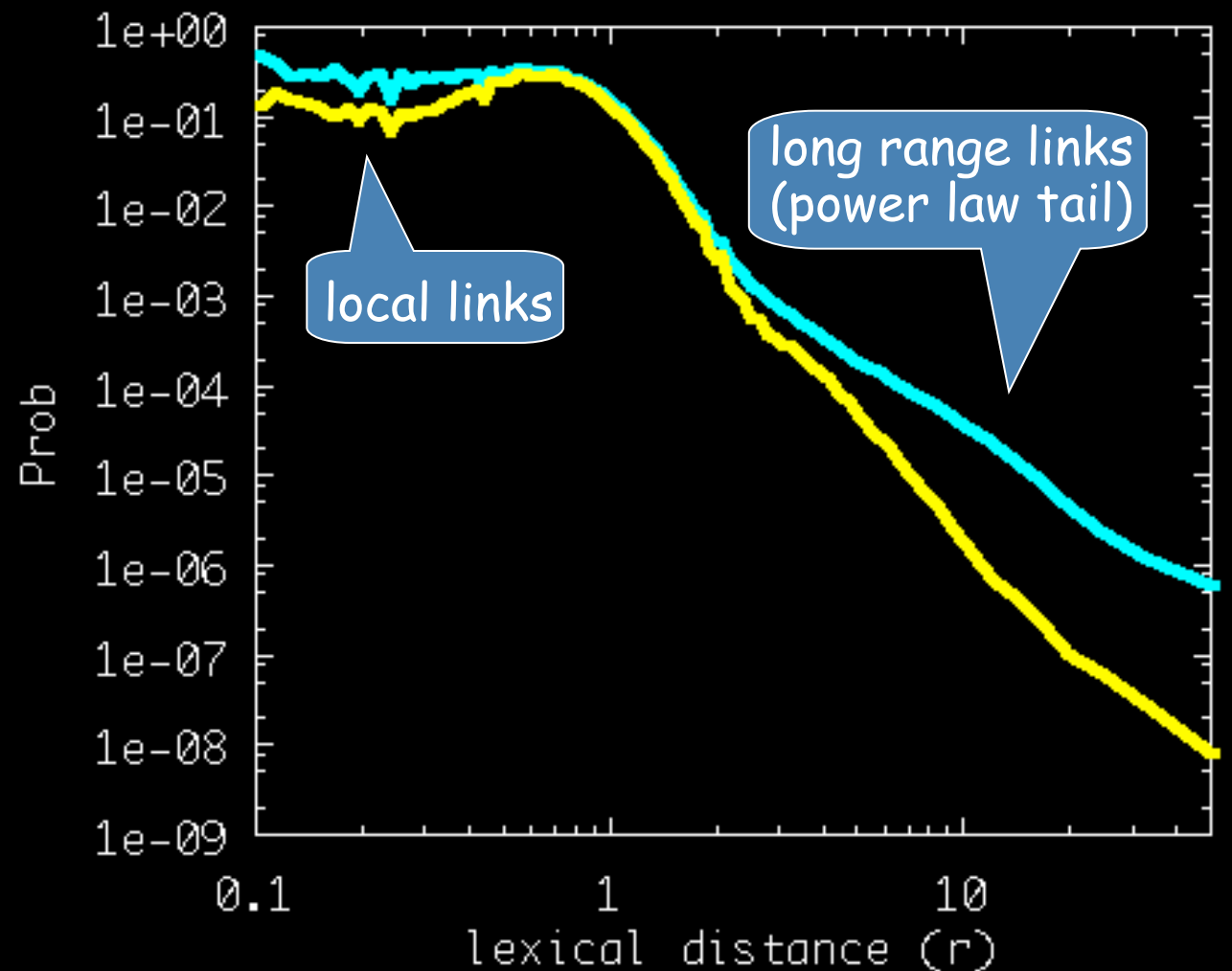
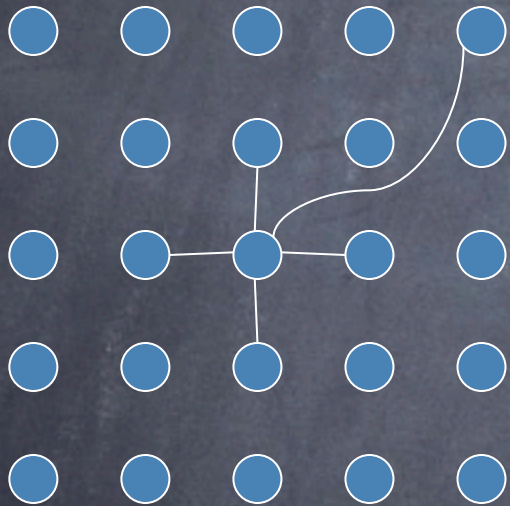
- Geographical networks (Kleinberg 2000)
 - Local links to all lattice neighbors
 - Long-range link probability distribution:
power law $P_r \sim r^{-\alpha}$
 - r : lattice (Manhattan) distance
 - α : constant clustering exponent

$$t \sim \log^2 N \Leftrightarrow \alpha = D$$

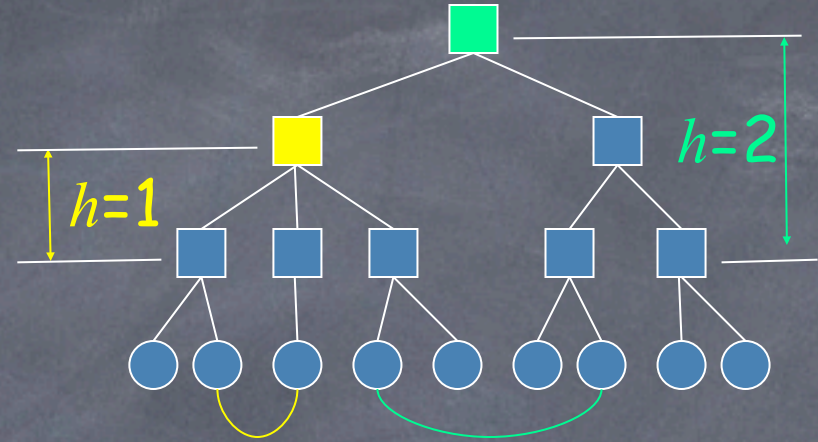
Is the Web a geographical network?

Replace lattice distance by lexical distance

$$r = (1 / \sigma_c) - 1$$



Exception # 2



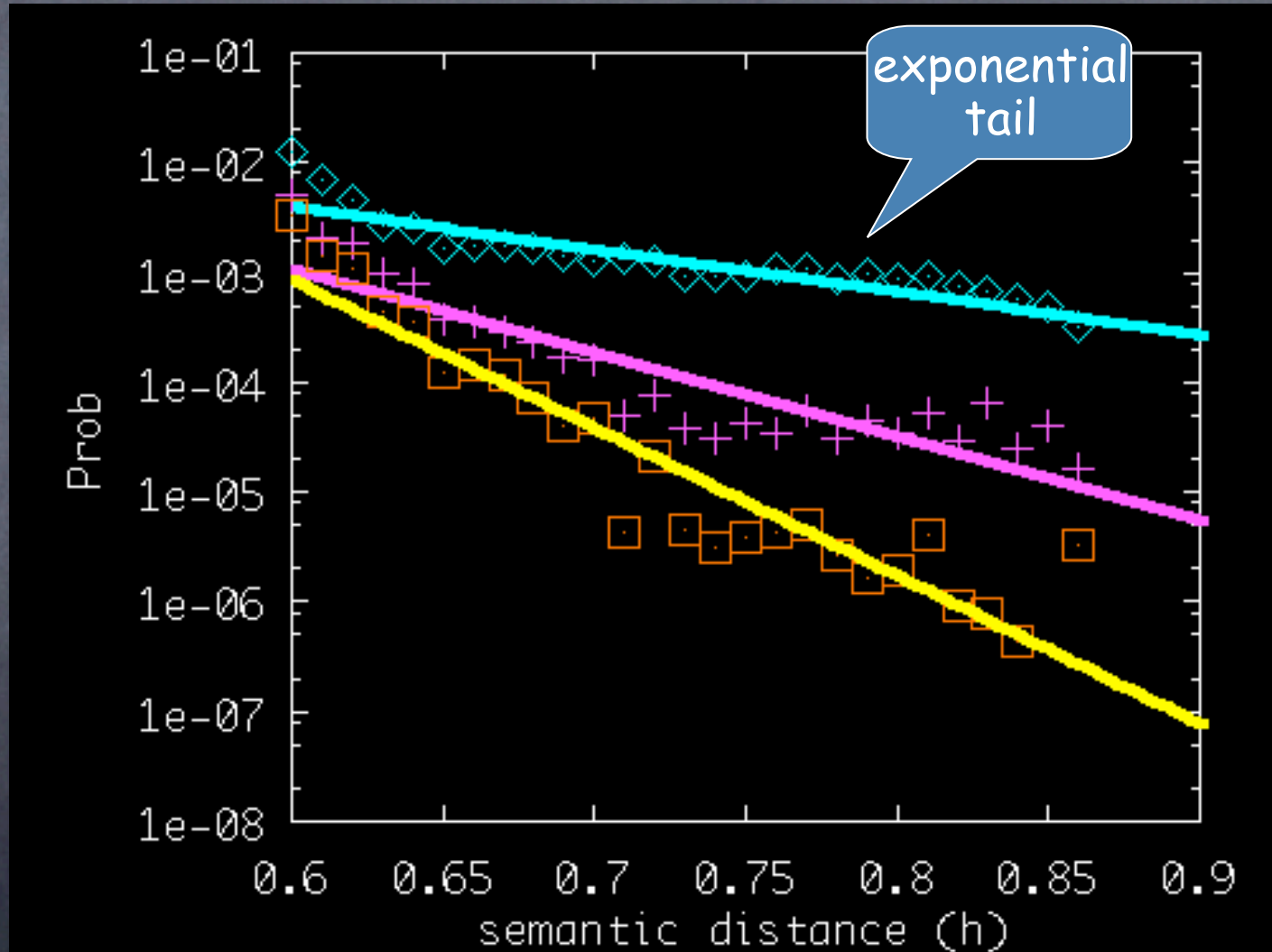
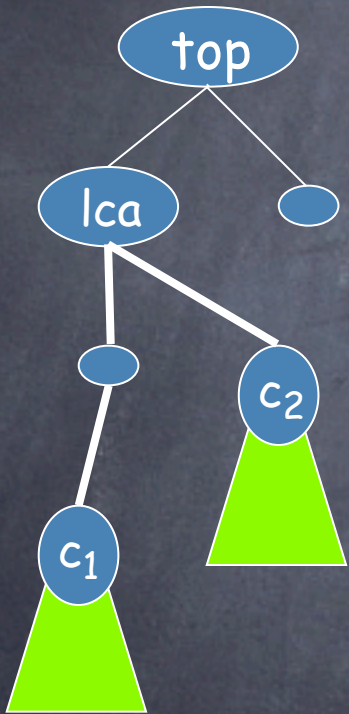
- Hierarchical networks
(Kleinberg 2002, Watts & al. 2002)
 - Nodes are classified at the leaves of tree
 - Link probability distribution: exponential tail
 $Pr \sim e^{-h}$
 - h : tree distance (height of lowest common ancestor)

$$t \sim \log^{\varepsilon} N, \varepsilon \geq 1$$

Is the Web a hierarchical network?

Replace tree distance by semantic distance

$$h = 1 - \sigma_S$$



Discussion: Does this
stuff really have any
applications?

Outline

- ✓ Mapping
 - > Topical locality
 - > Content, link, and semantic topologies in the Web
- ✓ Modeling
 - > How the Web evolves and why content matters
 - > Consequences for navigation and crawling
- ◉ Mining
 - > Topical Web crawlers
 - > Adaptive, intelligent crawling techniques
- ◉ Mingling
 - > Social Web search & recommendation
 - > Distributed collaborative peer search

Crawler applications

- **Universal Crawlers**
 - Search engines!
- **Topical crawlers**
 - Live search
(e.g., myspiders.informatics.indiana.edu)
 - Topical search engines & portals
 - Business intelligence (find competitors/partners)
 - Distributed, collaborative search

Topical crawlers

- Seminal work
 - Cho & Garcia-Molina @ Stanford
 - Chakrabarti & al @ IBM Almaden / IIT
 - Amento & al @ ATT Shannon
 - Ben-Shaul & al @ IBM Haifa
 - many others...

Topical crawlers

Miselanous Physics Websites

[sic]

- [Coaxial Cable Attenuation & Power Handling Calculator](#)
- [Britney Spears](#) guide to Semiconductor Physics: semiconductor physics, Edge Emitting Lasers and VCSELs
- [Particle Data Book](#) everything you ever wanted to know about particles, and then some.
- [X-Ray Data Booklet](#) everything you ever wanted to know about x-rays, and then some.

['A Brief History of Anglo-Saxon England'](#) -

['Anglo-Saxon Military Organisation'](#) - Article.

['Anglo-Saxon Social Organisation'](#) - Article.

['Arms and Armour - Part 3 Axes'](#) - Article.

['Arms and Armour - Part 7 Helmets'](#) - Article.

['Arms and Armour - Part 6 Mail Armour'](#) - Article.

['Arms and Armour - Part 4 Missile Weapons'](#) - Article.

['Arms and Armour - Part 2 Scramseaxes'](#) - Article.

['Arms and Armour - Part 8 Shields'](#) - Article.

['Arms and Armour - Part 1 Spears'](#) - Article.

['Arms and Armour - Part 5 Swords'](#) - Article.

['A Nice Little Earner'](#) - The slave trade in Anglo-Saxon England.

['A Spring Warmer'](#) - An alternative recipe for jugged hare!

['The Battle of Hastings'](#) - Article.

['Bone and Antler Working'](#) - Article.

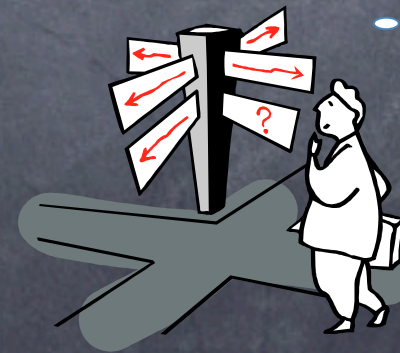
['Braid Weaving'](#) - Article.

['Bronzeworking'](#) - Article.

['Charcoal Burning'](#) - The results of an experiment in charcoal burning.

['Church Organisation'](#) - The organisation of the church in Anglo-Saxon England.

spears



Evaluating topical crawlers

- Goal: build “better” crawlers to support applications
- Build an unbiased **evaluation framework**
 - Define **common tasks** of measurable difficulty
 - Identify **topics**, relevant **targets**
 - Identify appropriate performance measures
 - **Effectiveness**: quality of crawler pages, order, etc.
 - **Efficiency**: separate CPU & memory of crawler algorithms from bandwidth & common utilities

Evaluating topical crawlers: Topics

Keywords

Description

Targets

dmoz open directory project

Home: [Cooking](#): [Baking and Confections](#): [Cookies](#): [Chocolate Chip](#) (6)

- [The Big Chocolate Chip Cookie Page](#) - Devoted to the chocolate chip cookie.
- [Chocolate Chip Cookies](#) - Various recipes for cookies with morsels of chocolate.
- [Chocolate Chip Cookies from Allrecipes](#) - Include regular, nuts, white chocolate.
- [In the Chips](#) - Cookies, cakes, candy, muffins, etc. using chocolate chips.

Copyright © 1998-2001 Netscape

[Terms of Use](#)

- Automate evaluation using edited directories
- Different sources of relevance assessments

Recipe and Tips R

A chocolate chip cookie recipe that uses Karo syrup in it.
 A recipe using metric measurements?
 Any recipes that don't use eggs?
 A chocolate chip cookie recipe that you bake in mini muffin pans w
 How does one avoid dry, 'cakey' cookies?
 Any recipes for Chocolate Chip Coolie Pies?
 The recipe for chocolate chip cookies in a jar. All the dry ingredien

Recipe	Rating
Absolutely Excellent Oatmeal Cookies Submitted by: Marylou	★★★★★ 46 Ratings 25 Reviews
Absolutely Sinful Chocolate Chocolate Chip Cookies Submitted by: Marsha	★★★★★ 63 Ratings 49 Reviews

Cookies that are out of this world...



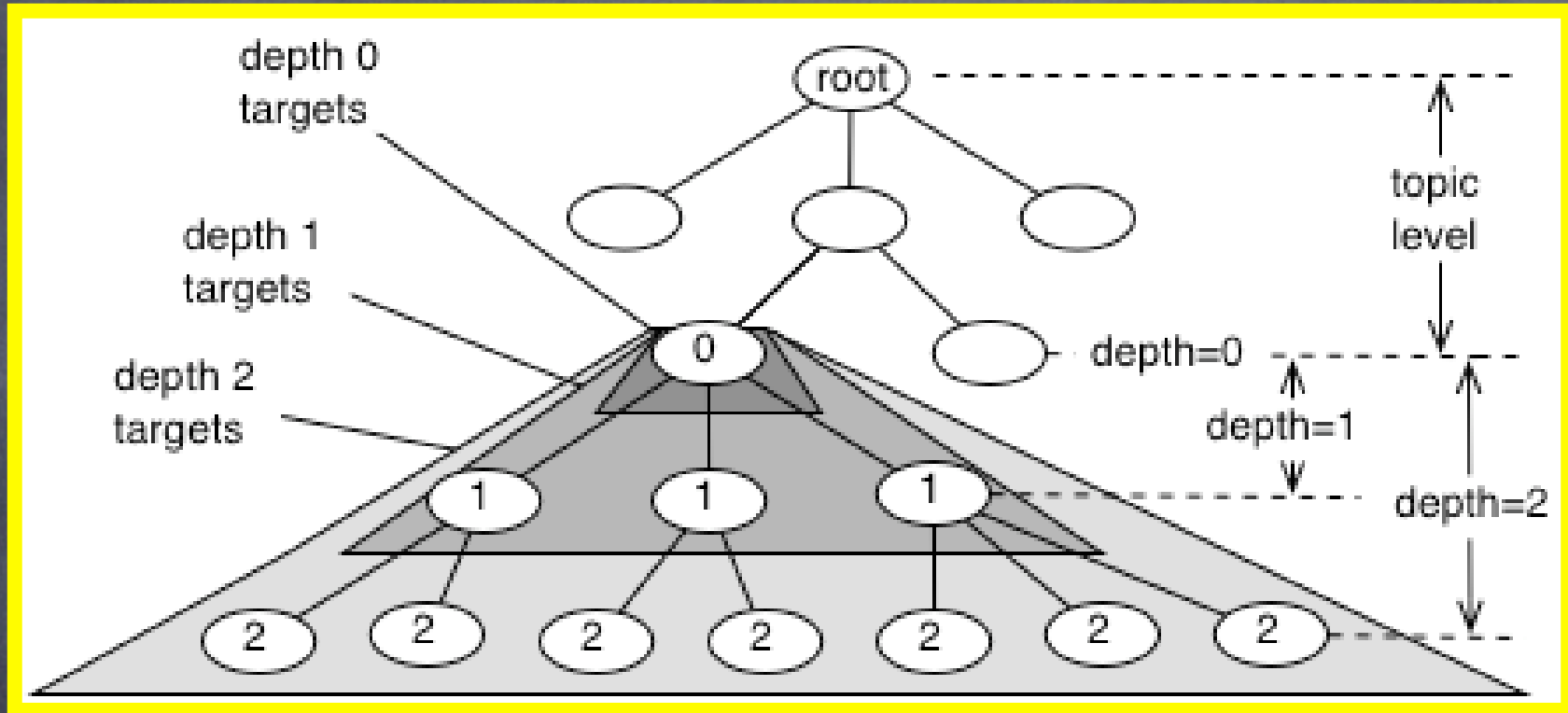
In the kitchen of a Whitman Massachusetts country inn, the first chocolate chip cookie emerged in 1937. Simple experiments led to a recipe combining bits of chocolate candy with a kind of butter cookie cookie dough resulting in a delicious mixture that offered the crunchiness of a cookie with a taste of chocolate candy in every bite. Obviously the cookies were a hit at the Inn and wherever else the recipe spread. Chocolate chip cookies have remained an American homemade treat.

CHOCOLATE CHIP COOKIES

RECIPE INDEX

- [BLACK AND WHITE CHOCOLATE CHIPPERS](#)
- [CLASSIC CHOCOLATE CHIP COOKIES](#)
- [COW CHIP COOKIES](#)
- [DEVIL'S FOOD CHOCOLATE CHIP COOKIES](#)
- [GOTTA HAVE EM' NOW! COOKIES](#)
- [MINT CHOCOLATE SANDWICH COOKIES](#)
- [NEIMAN MARCUS CHOCOLATE CHIP COOKIES](#)
- [OLD FASHIONED CHOCOLATE CHIPPERS](#)

Topics and Targets

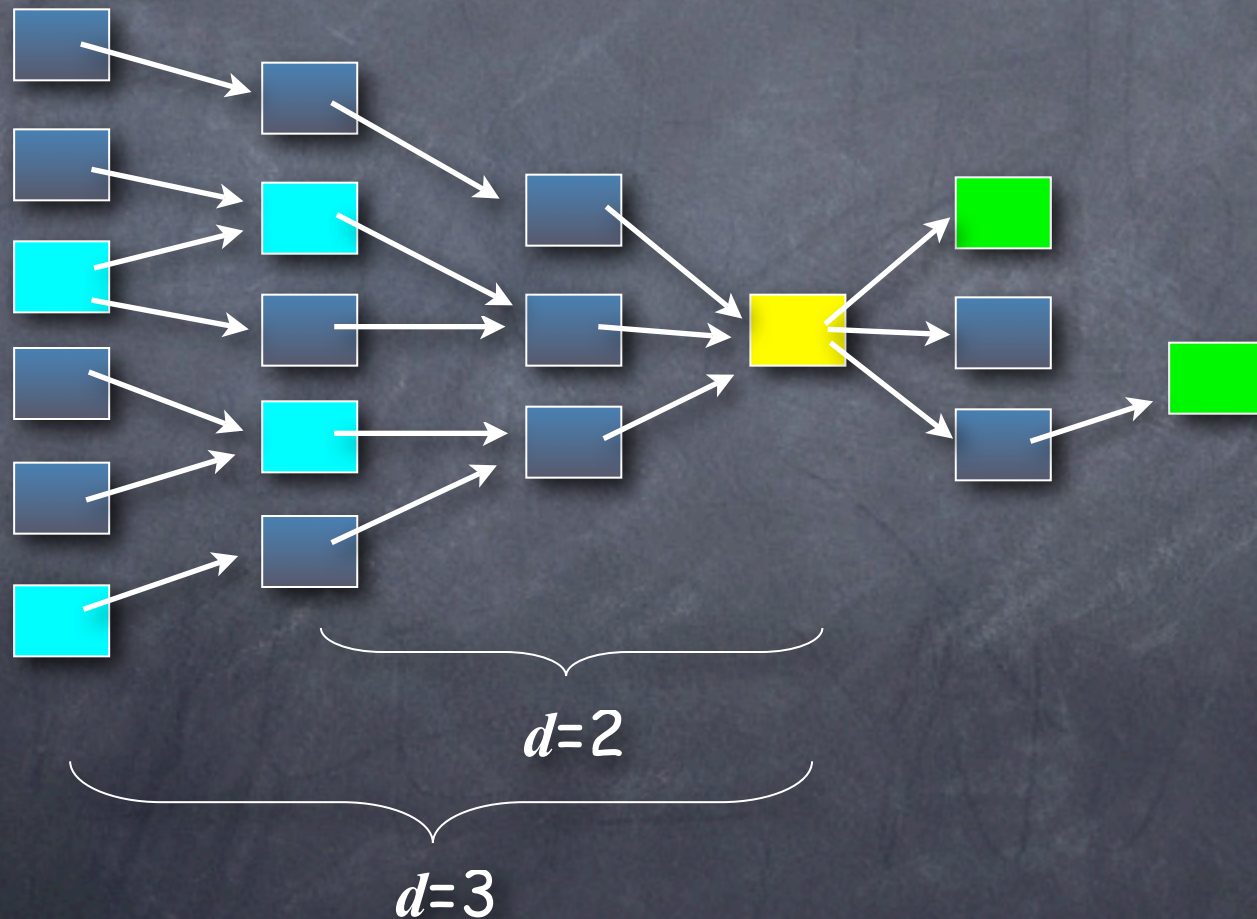


topic level ~ specificity

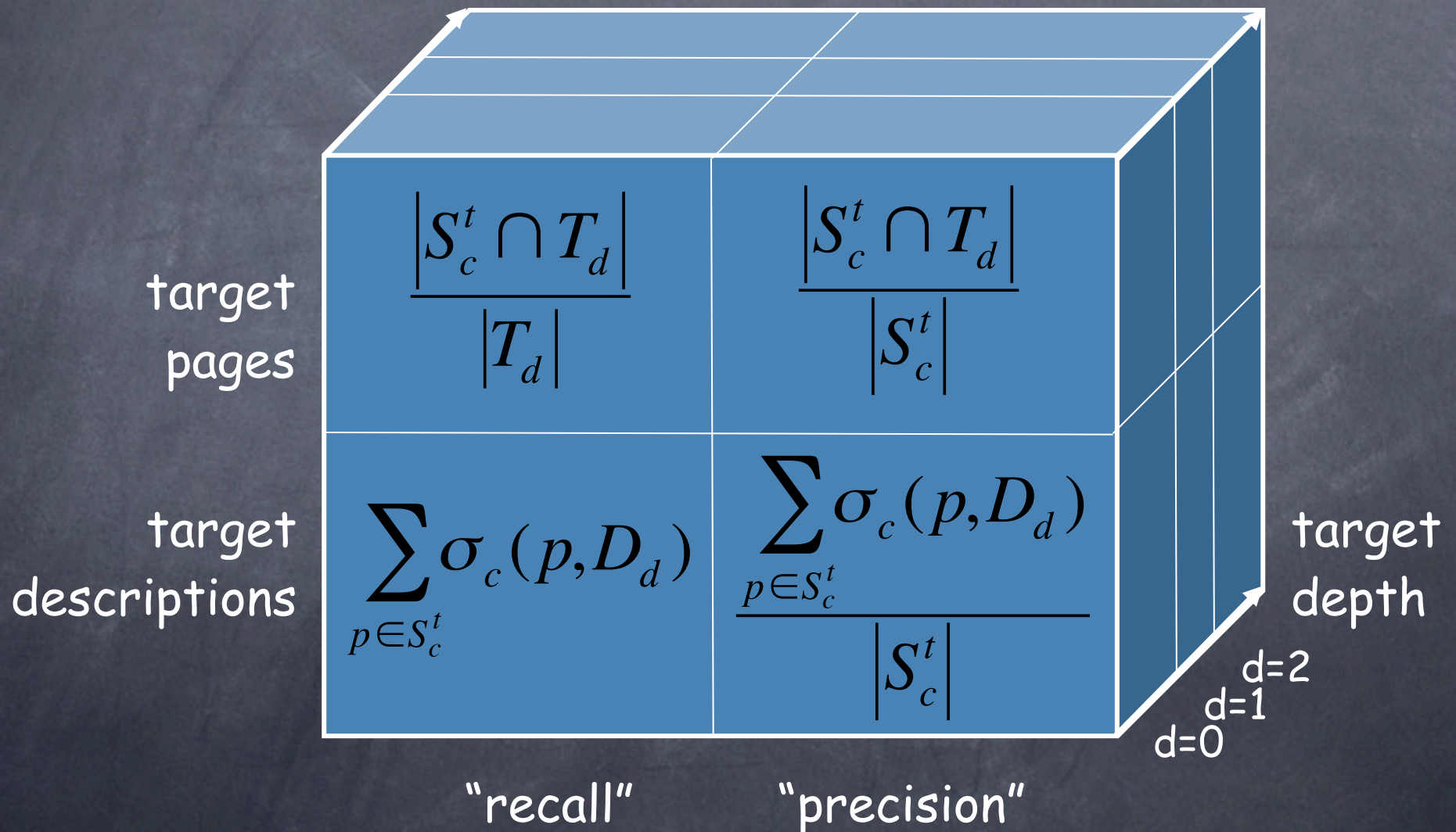
depth ~ generality

Evaluating topical crawlers: Tasks

Start from **seeds**, find **targets**
and/or pages **similar to target descriptions**

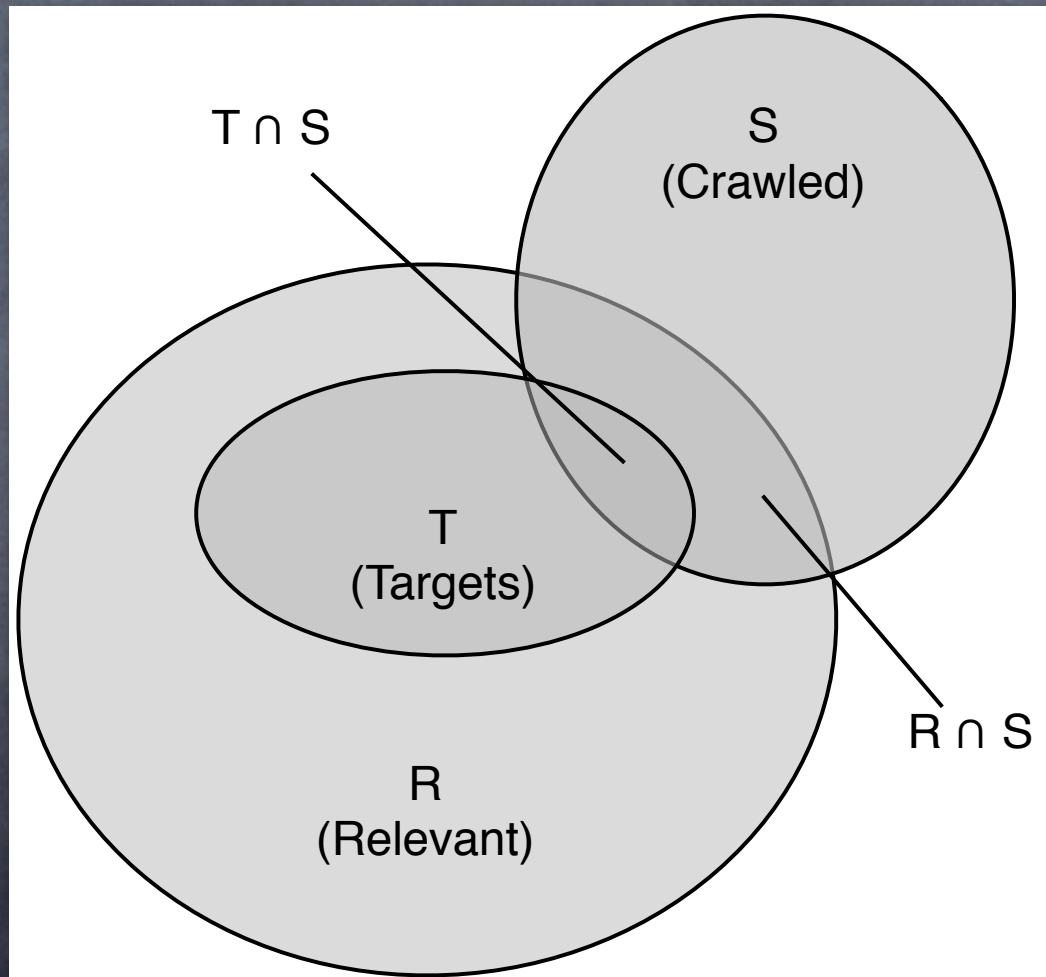


Performance matrix

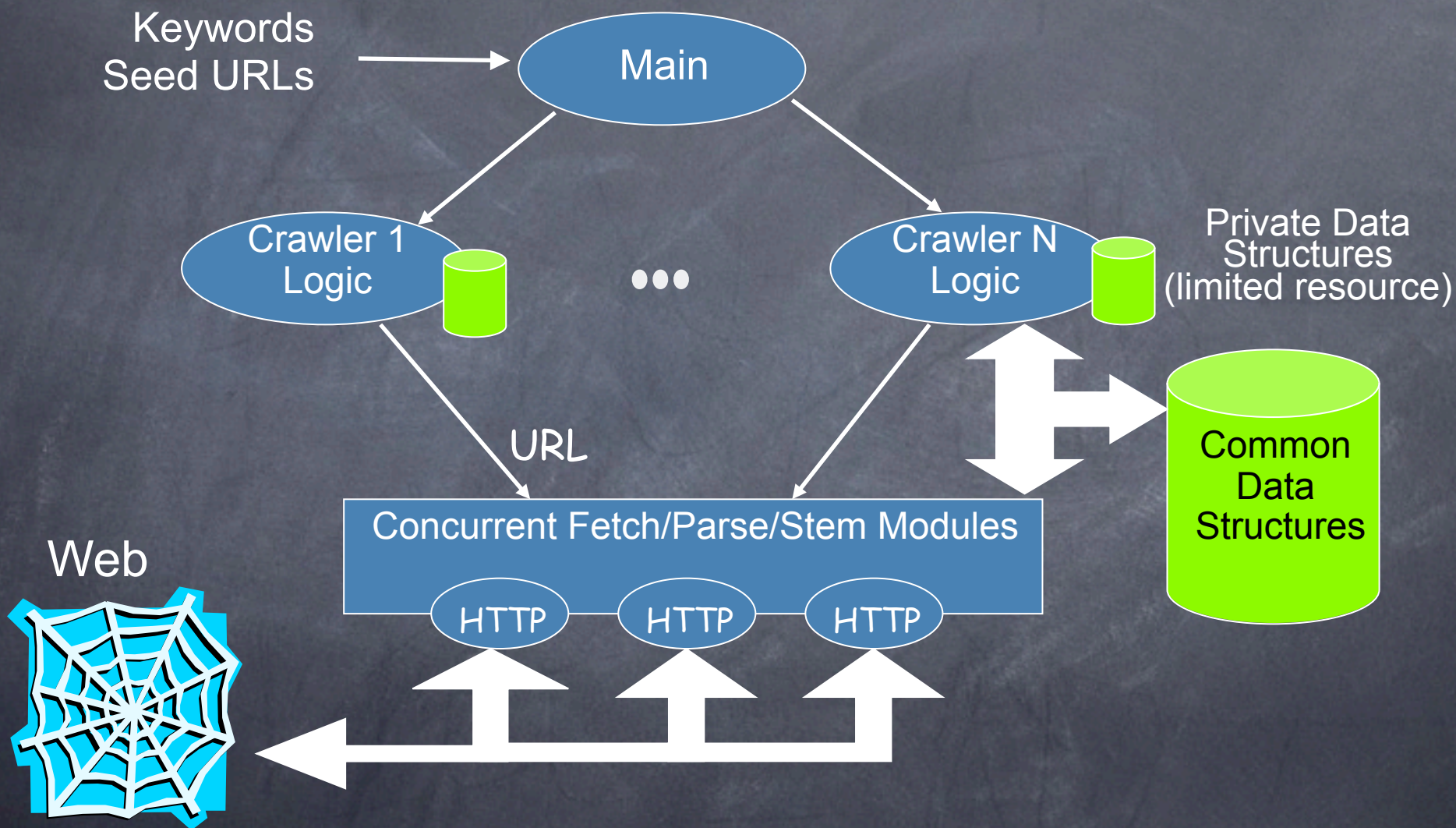


Discussion:

What assumption (targets)?



Evaluating topical crawlers: Architecture

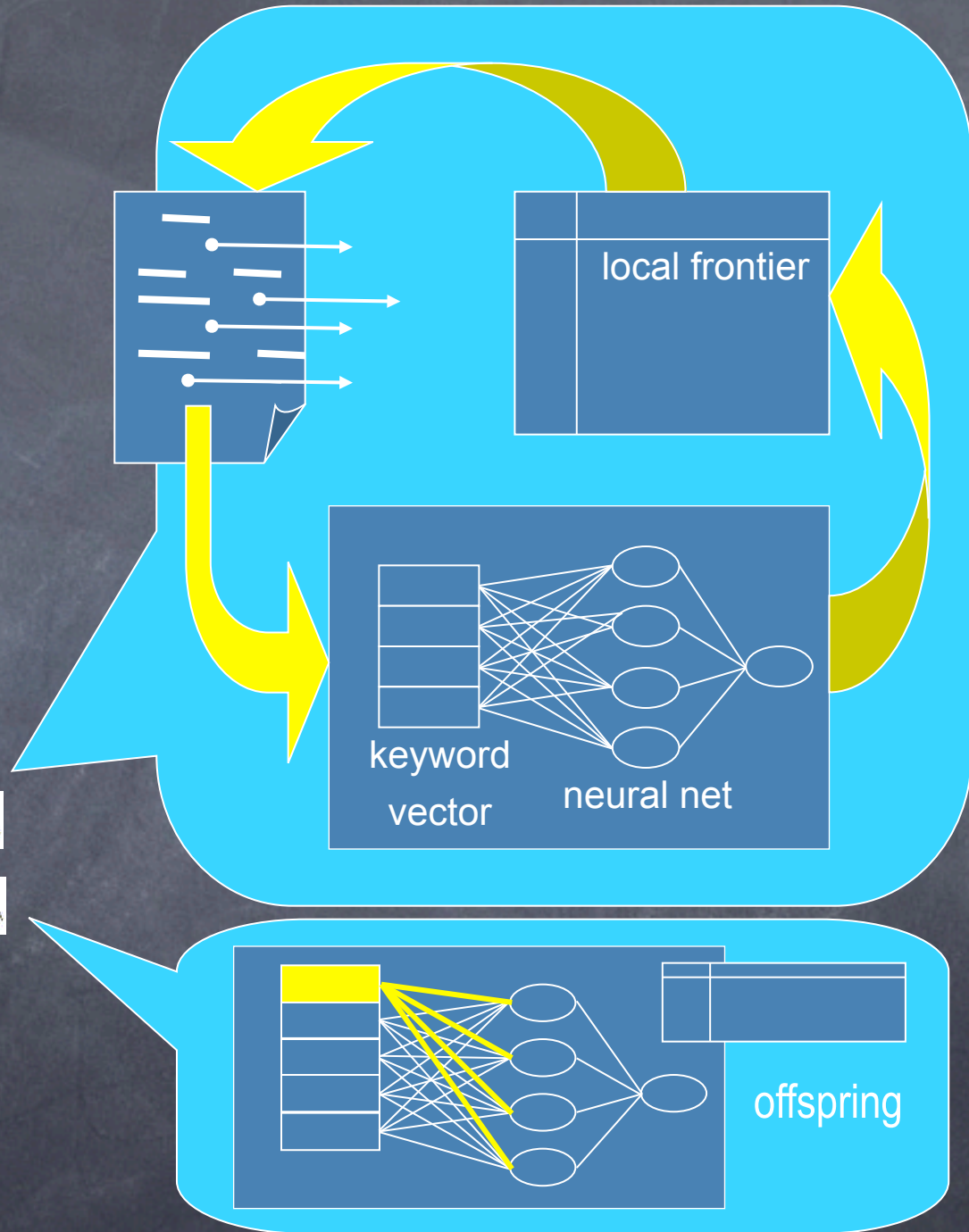
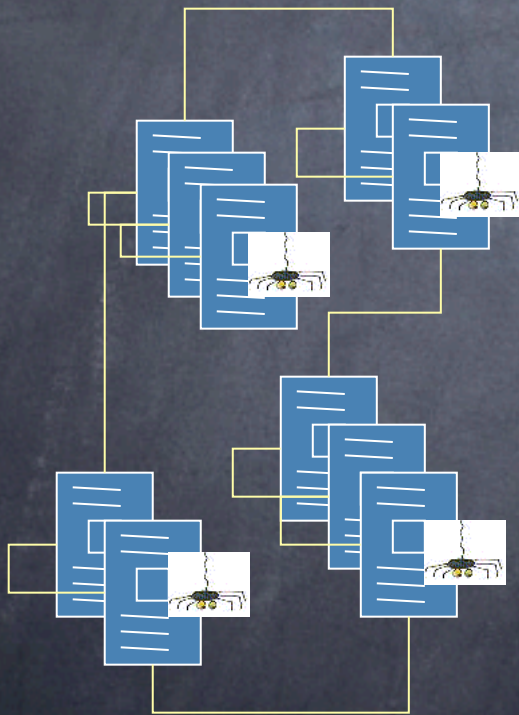


Examples of crawling algorithms

- **Breadth-First**
 - Visit links in order encountered
- **Best-First**
 - Priority queue sorted by similarity
 - Variants:
 - explore top N at a time
 - tag tree context
 - hub scores
- **SharkSearch**
 - Priority queue sorted by combination of similarity, anchor text, similarity of parent, etc.
- **InfoSpiders**

InfoSpiders

adaptive distributed
algorithm using an
evolving population of
learning agents



InfoSpiders



Author.....FILIPPO MENCZER

Advisor.....RIK BELEW

University of California, San Diego

Easy query:
The population rapidly focuses
on the relevant areas of the
information space

InfoSpiders



Author.....FILIPPO MENCZER

Advisor.....RIK BELEW

University of California, San Diego

**Impossible query:
The population goes extinct**

InfoSpiders



Author.....FILIPPO MENCZER

Advisor.....RIK BELEW

University of California, San Diego

**Ambiguous query:
A subpopulation eventually
locates the relevant pages**

Mozilla Firefox

http://myspiders.informatics.indiana.edu/myspiders2.html

Query: search censorship in france

Start Stop Max. Pages: 100

MySpiders

Crawler Name: InfoSpiders Pages Crawled: 100

Population: 0

Source	URL	Score	Rece...
Spider2	http://www.multilingual-search.com/new-tool-shows-cens...	0.43	?
Seed	http://www.technologynewsdaily.com/node/2283	0.34	1
Seed	http://www.rpi.edu/~bulloj/search/CENSORSHIP.html	0.33	0
Seed	http://en.wikipedia.org/wiki/Censorship_in_France	0.32	0.14
Seed	http://www.laboratorytalk.com/news/iqd/iqd109.html	0.3	0.33
Seed	http://www.informatics.indiana.edu/news/news.asp?id=313	0.27	?
Seed	http://blog.searchenginewatch.com/blog/050117-090638	0.12	0.02

Spider Hierarchy

- Spiders
 - Spider1
 - Spider2
 - Spider13
 - Spider3
 - Spider4
 - Spider5
 - Spider11
 - Spider12
 - Spider6
 - Spider7

New tool shows censorship by search engine in China, France, Germany and the US [Multilingual Search]

New tool shows censorship by search engine in China, France, Germany and the US



Andy Atkins-Krüger Mar 19, 2006 | [en]

Pandia reports on a new censorship comparison tool developed by researchers at Indiana University. The tool compares the preeminence of words featured in the top ten results of Yahoo or Google by displaying words graphically giving weight to those terms which are more frequent - the end result is a graphic somewhat reminiscent of Technorati tags.

CENSEARCHIP

Called **Censearchip** - the team behind the tool have clearly chosen China, France and Germany - compared with the US - because of the recent censorship issues in China and the restrictions placed on search engines in terms of displaying nazi material - by Germany and France.

In addition to this 'political' censorship - it would be useful to be able to compare results by country where 'business' censorship has an impact - for instance the filtering of results which takes place between the US and the UK.

Spider Details

- Details
 - Spider13
 - Status
 - Energy
 - 1.1394327
 - Query
 - Term1
 - Term2
 - Term3
 - Term4
 - itali
 - History

Spider Details

- Details
 - Spider11
 - Status
 - Energy
 - Query
 - Term1
 - Term2
 - Term3
 - Term4
 - engin
 - History

Evolutionary Local Selection Algorithm (ELSA)

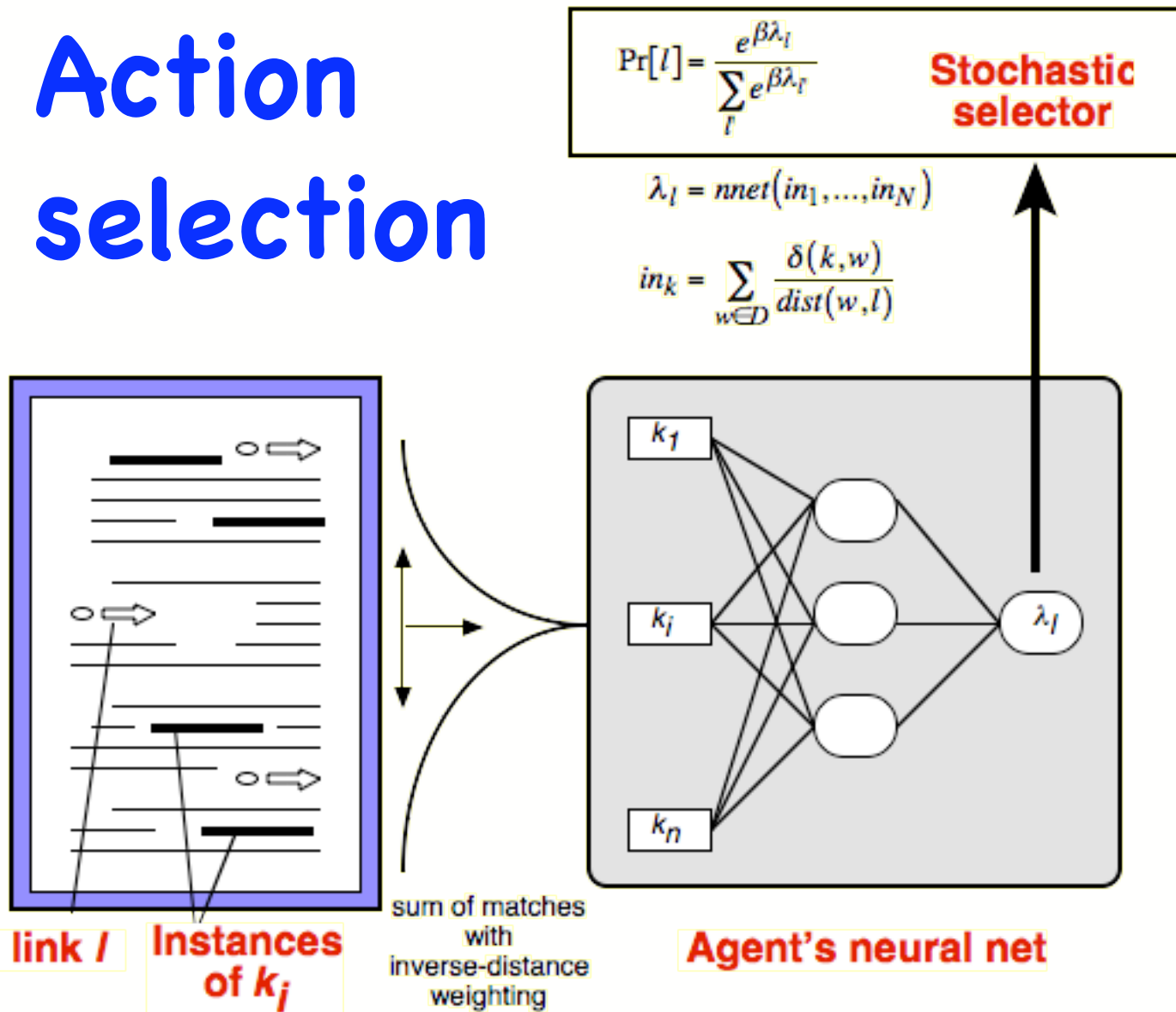
```
Foreach agent thread:  
  Pick & follow link from local frontier  
  Evaluate new links, merge frontier  
  Adjust link estimator  
   $E := E + \text{payoff} - \text{cost}$   
  If  $E < 0$ :  
    Die  
  Elsif  $E > \text{Selection\_Threshold}$ :  
    Clone offspring  
    Split energy with offspring  
    Split frontier with offspring  
    Mutate offspring
```

reinforcement
learning

match
resource
bias

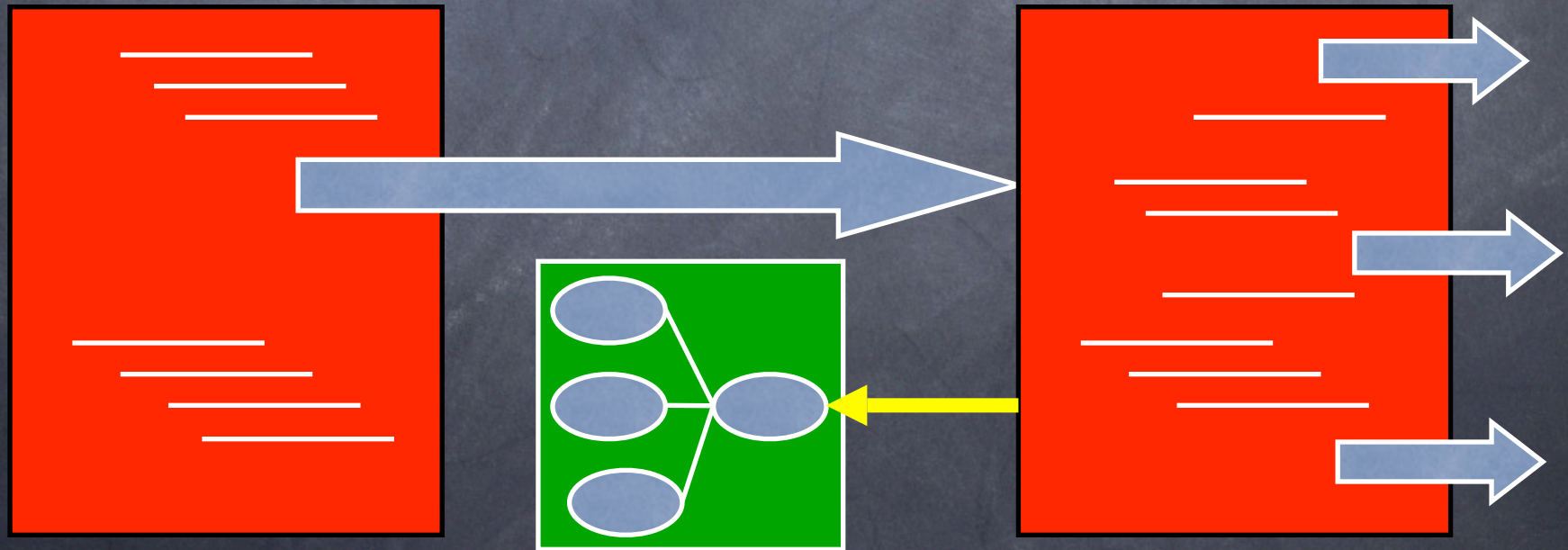
selective
query
expansion

Action selection

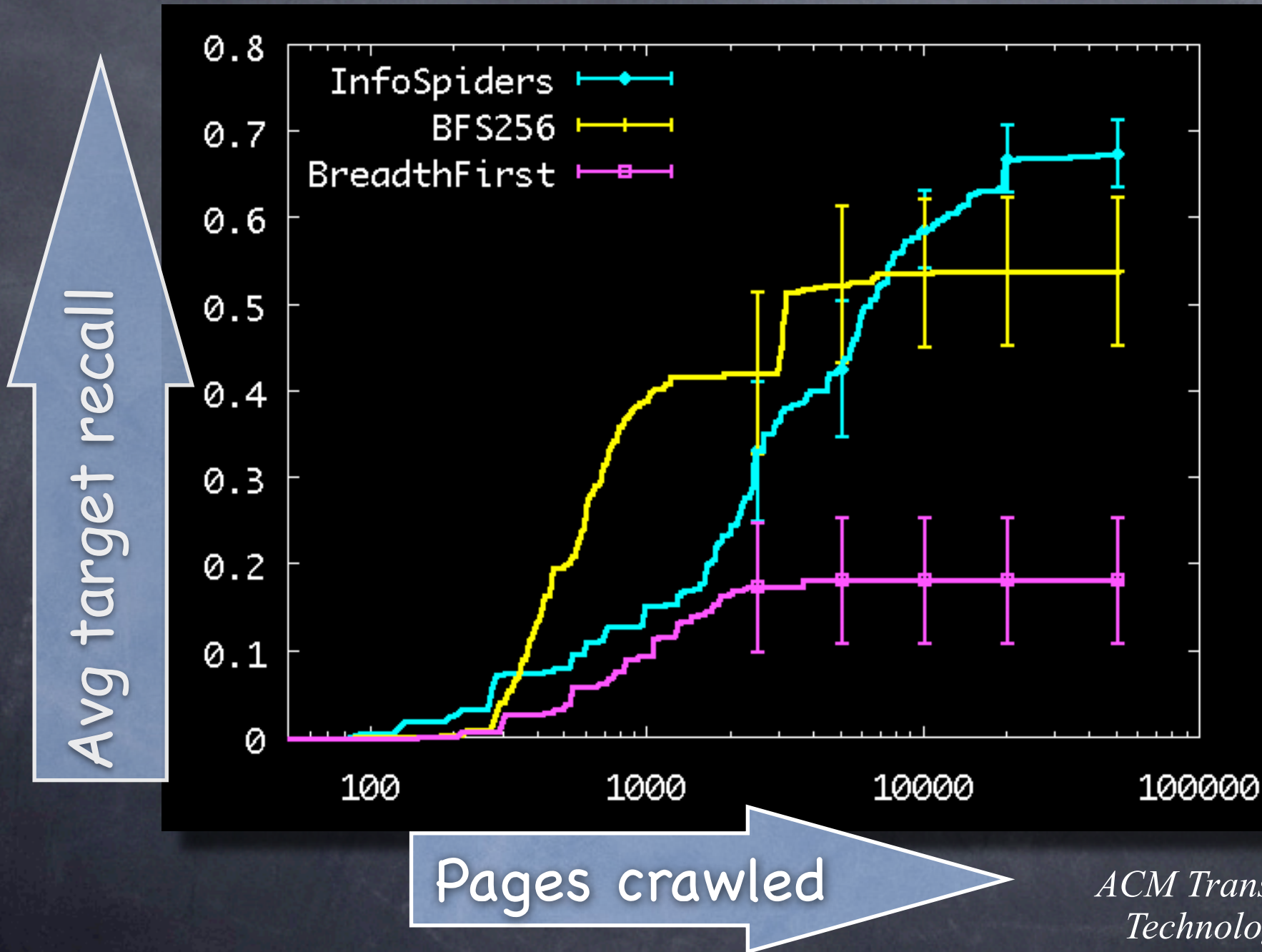


Q-learning

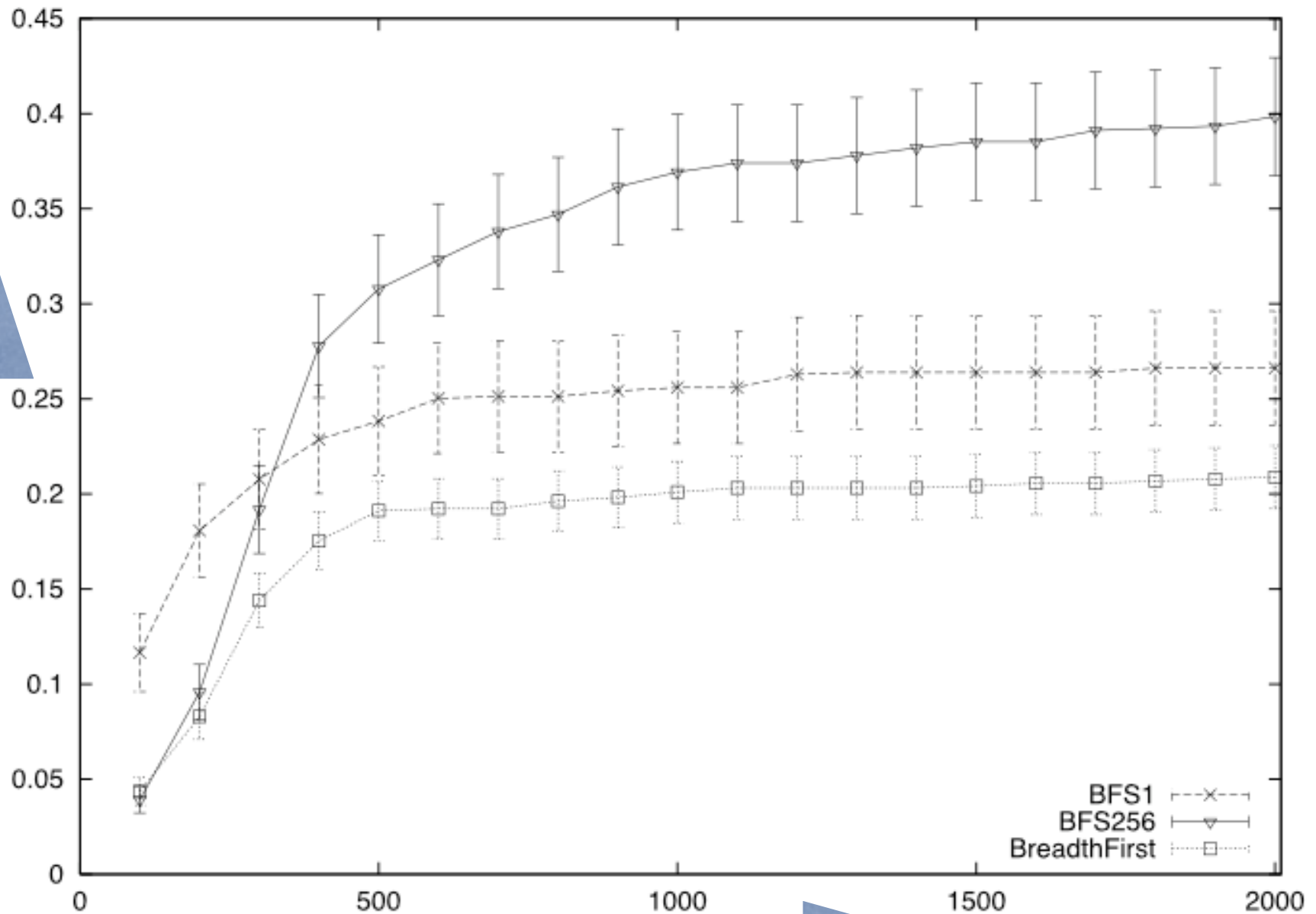
- Compare estimated relevance of visited document with estimated relevance of link followed from previous page
- Teaching input: $E(D) + \mu \max_{I(D)} \lambda_I$



Performance



Exploration vs. Exploitation

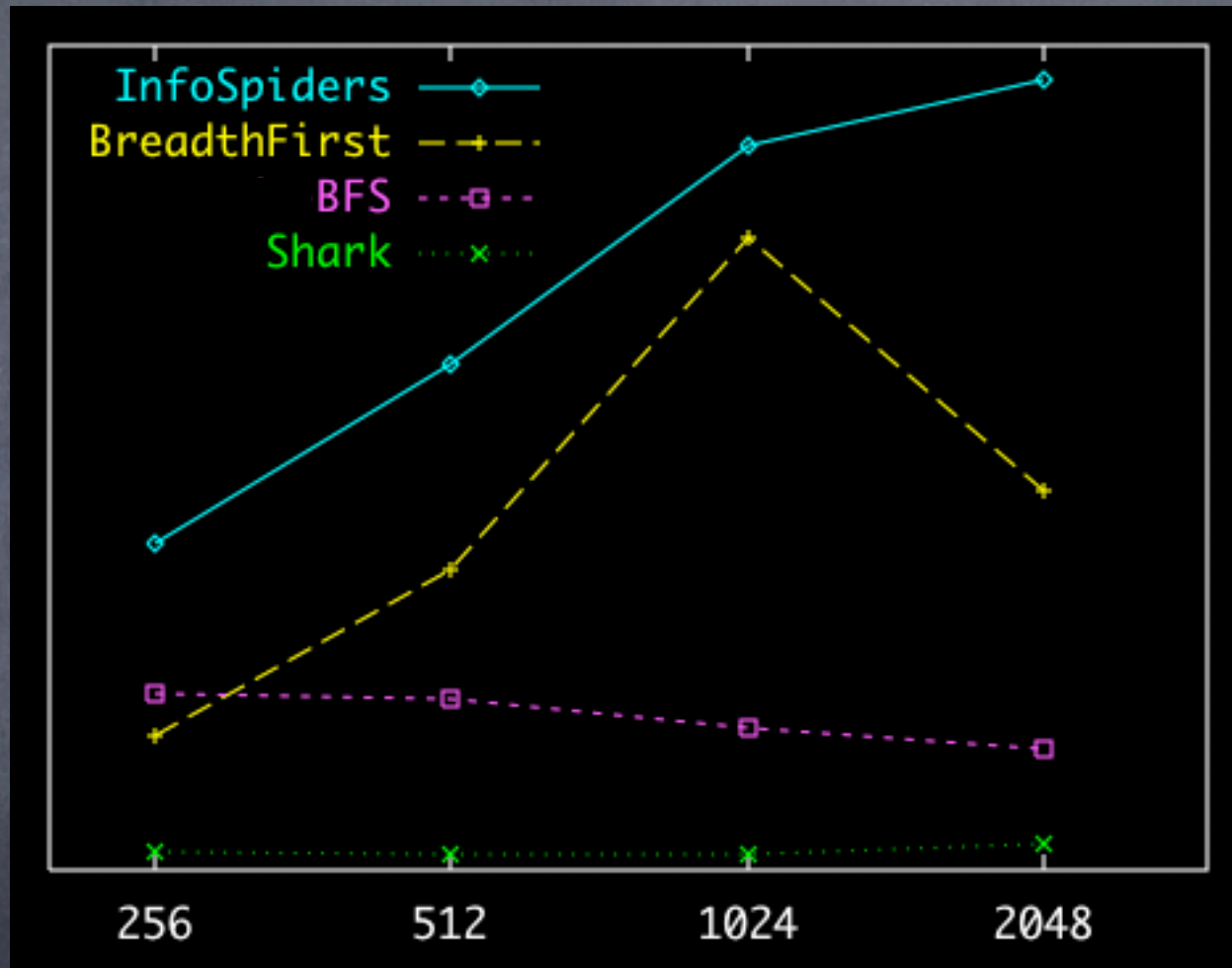


Avg target recall

Pages crawled

Efficiency & scalability

Performance/cost



Link frontier size

About Exelixis

Exelixis, Inc. is a leading genomics-based drug discovery company focused on product development through its expertise in comparative genomics and model system genetics. These technologies provide a rapid, efficient and cost effective way to move from DNA sequence data to knowledge about the function of genes and the proteins they encode. The company's technology is broadly applicable to all life sciences industries including pharmaceutical, diagnostic, agricultural biotechnology and animal health. Exelixis has partnerships with Aventis CropScience S.A., Bayer Corporation, Bristol-Myers Squibb Company, Elan Pharmaceuticals, Inc., Pharmacia Corporation, Protein Design Labs, Inc., Scios Inc. and Dow AgroSciences LLC, and is building its internal development program in the area of oncology. For more information, please visit the company's web site at www.exelixis.com.



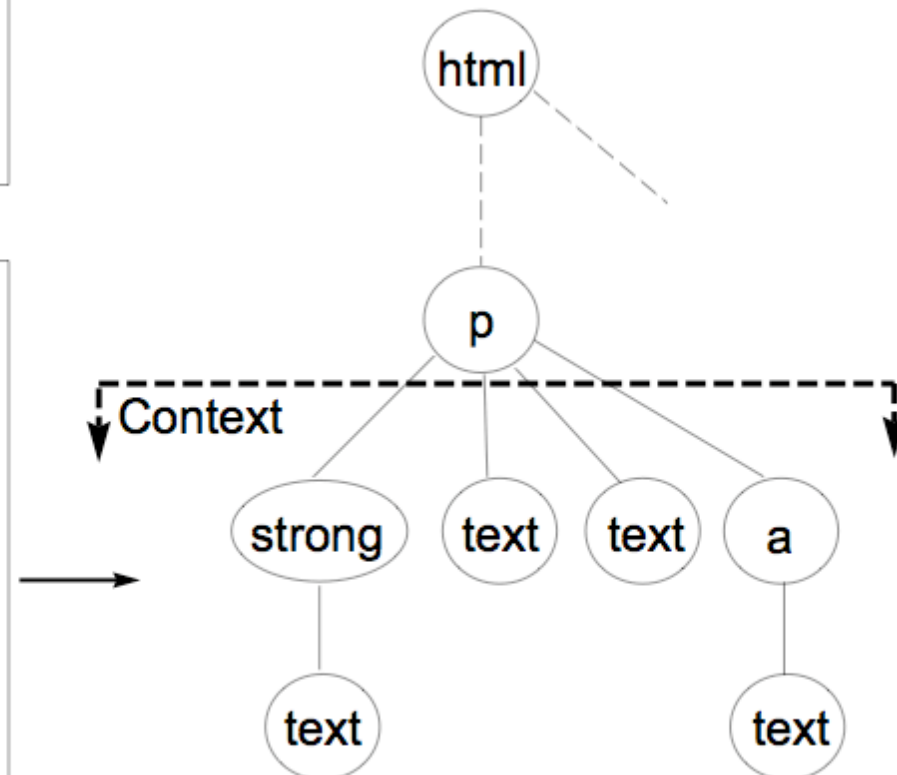
```
<P class=MsoNormal>
<STRONG>About Exelixis</STRONG><BR>Exelixis, Inc. is a leading genomics-based drug discovery
company focused on product development through its expertise in comparative genomics and model system
genetics. These technologies provide a rapid, efficient and cost effective way to move from DNA sequence
data to knowledge about the function of genes and the proteins they encode. The company's technology is
broadly applicable to all life sciences industries including pharmaceutical, diagnostic, agricultural
biotechnology and animal health. Exelixis has partnerships with Aventis CropScience S.A., Bayer Corporation,
Bristol-Myers Squibb Company, Elan Pharmaceuticals, Inc., Pharmacia Corporation, Protein Design Labs,
Inc., Scios Inc. and Dow AgroSciences LLC, and is building its internal development program in the area of
oncology.<SPAN style=mso-spacerun: yes>&nbsp;</SPAN></SPAN>For more information, please visit the
company's web site at
<A href="http://www.exelixis.com/">www.exelixis.com</A>.<o:p></o:p>
</P>
```



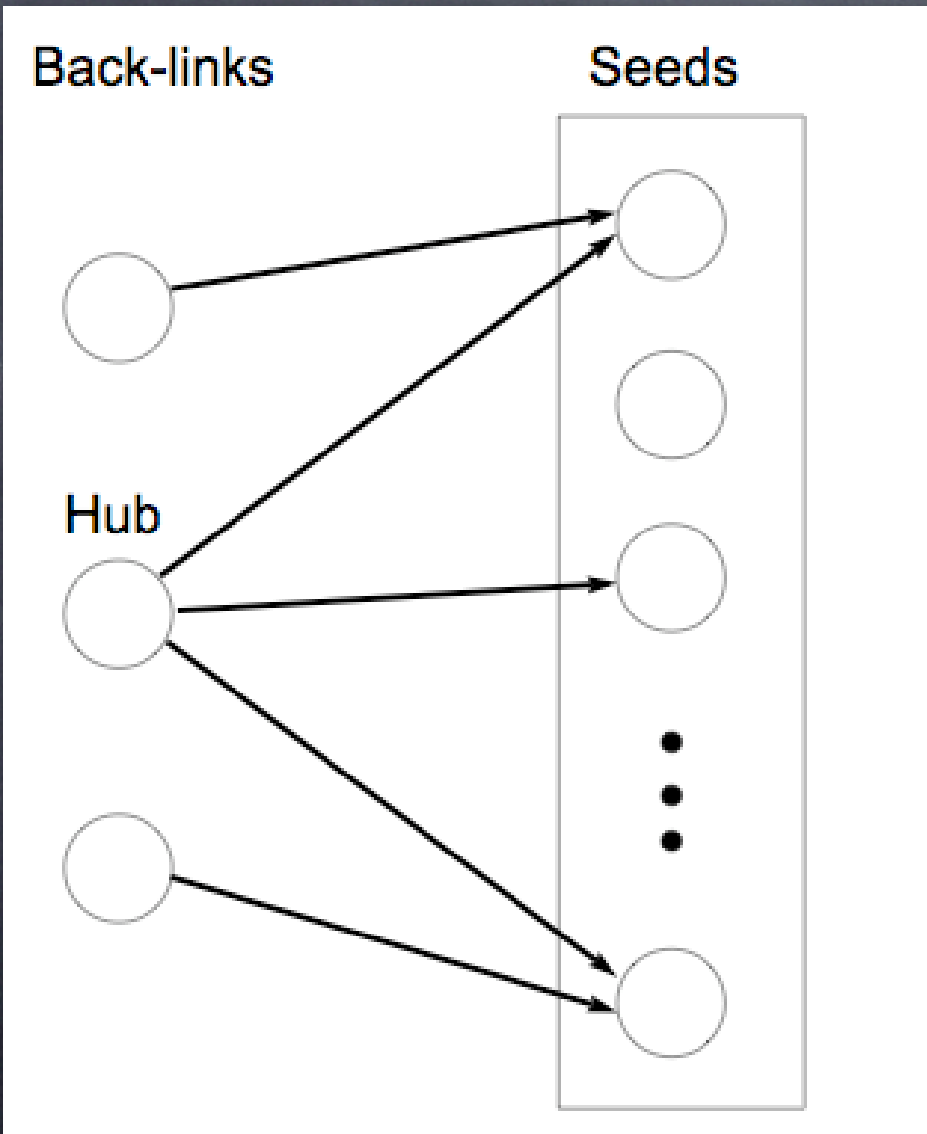
```
- <p>
- <strong>
  <text>about exelixis</text>
</strong>
<text>exelixis inc is a leading genomics based drug discovery company focused on
product development through its expertise in comparative genomics and model
system genetics these technologies provide a rapid efficient and cost effective
way to move from dna sequence data to knowledge about the function of genes
andthe proteins they encode the company s technology is broadly applicable to
all life sciences industries including pharmaceutical diagnostic agricultural
biotechnology and animal health exelixis has partnerships with aventis
cropscience s a bayer corporation bristol myers squibb company elan
pharmaceuticals inc pharmacia corporation protein design labs inc scios inc and
dow agrosiences llc and is building its internal development program in the area
of oncology</text>
<text>for more information please visit the company s web site at</text>
- <a href="http://www.exelixis.com/">
  <text>www exelixis com</text>
</a>
</p>
```

DOM context

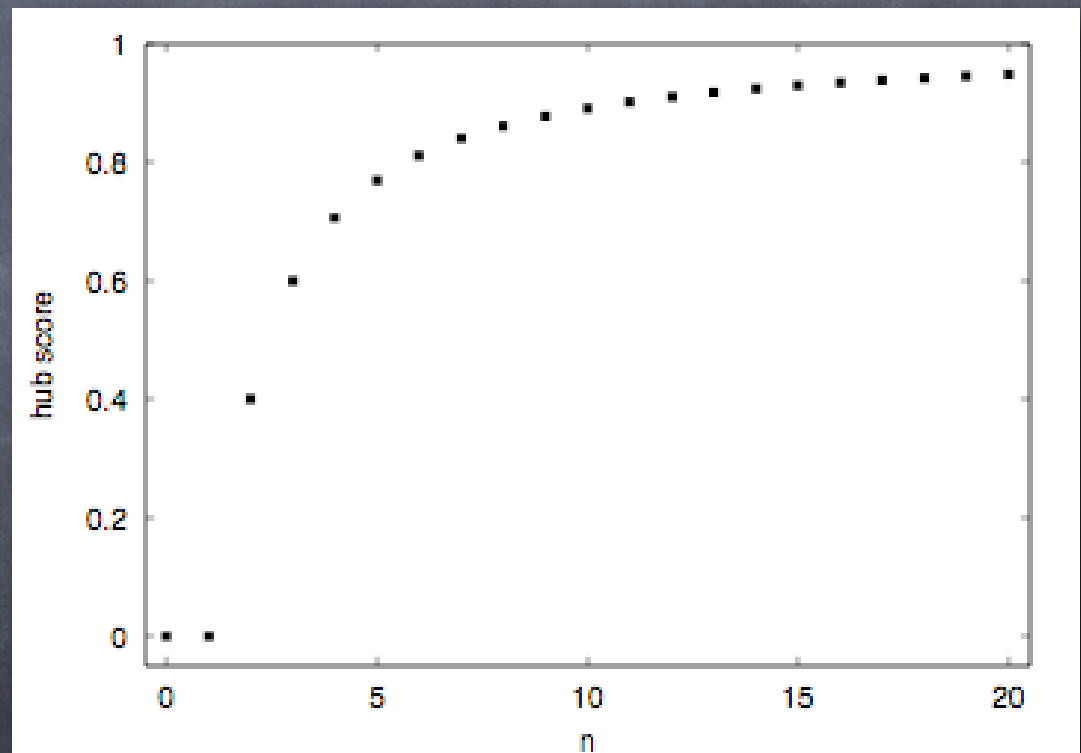
Link score = linear
combination between
page-based and
context-based
similarity score



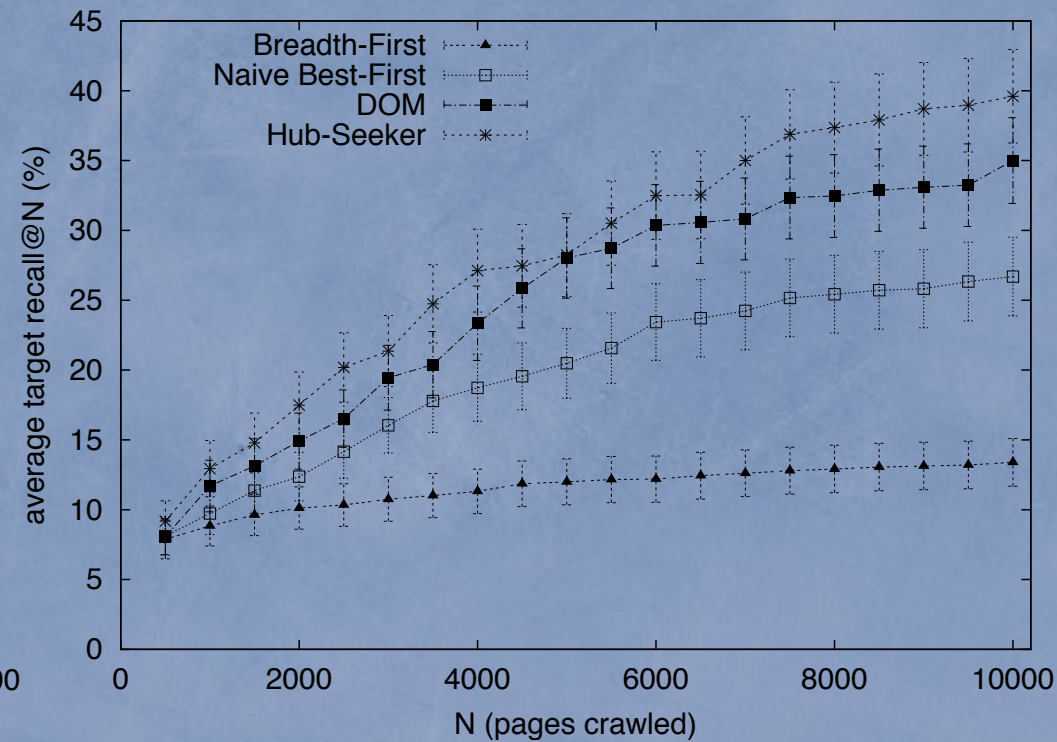
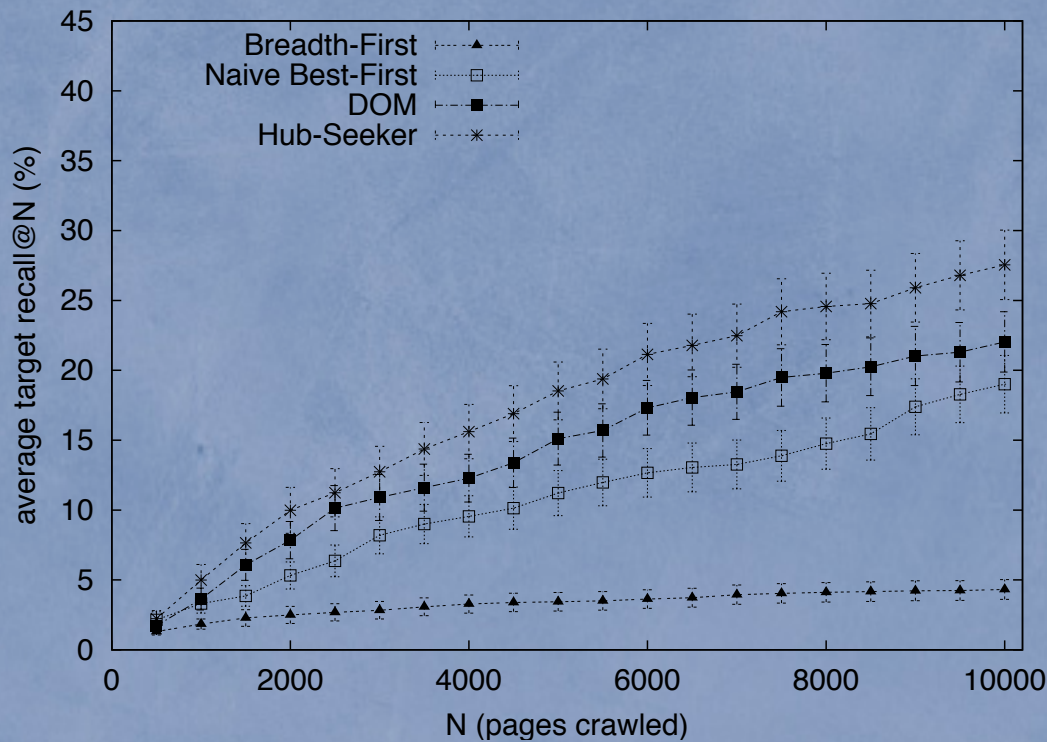
Co-citation: hub scores



Link score_{hub} = linear combination between link and hub score



Recall (159 ODP topics)



Split ODP URLs
between seeds and
targets

Add 10 best hubs to
seeds for 94 topics

Normative lessons about crawlers

- **Remember** links from previous pages
- **Exploration** is important
- **Adaptation** can help
 - to focus search
 - to expand query
 - to learn link estimates
- **Distributed** algorithms
 - boost efficiency
 - but watch for premature exploitation
- Improve **link evaluation** by looking at
 - parent pages
 - DOM context

Discussion: Which crawler algorithm would you use in your app?

Outline

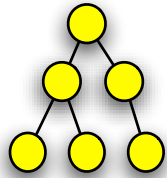
- ✓ Mapping
 - > Topical locality
 - > Content, link, and semantic topologies in the Web
- ✓ Modeling
 - > How the Web evolves and why content matters
 - > Consequences for navigation and crawling
- ✓ Mining
 - > Topical Web crawlers
 - > Adaptive, intelligent crawling techniques
- ◉ Mingling
 - > Social Web search & recommendation
 - > Distributed collaborative peer search



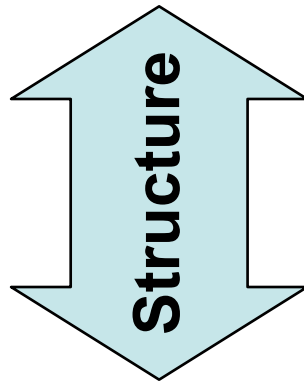
givealink.org



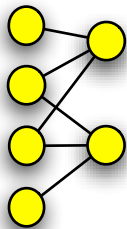
76788 links
1015 donations
117 registered users
Last updated: Apr 17, 2006



hierarchical

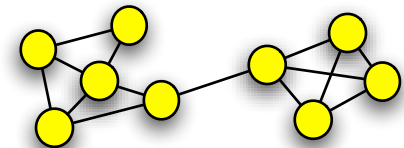
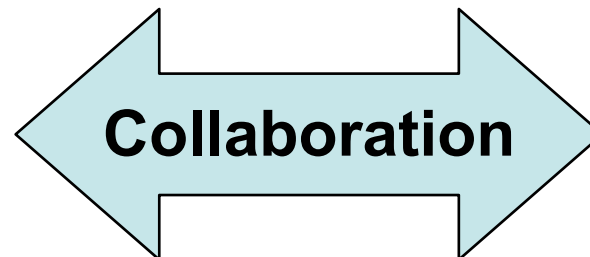
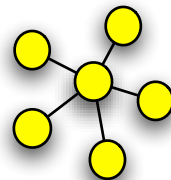


flat



centralized

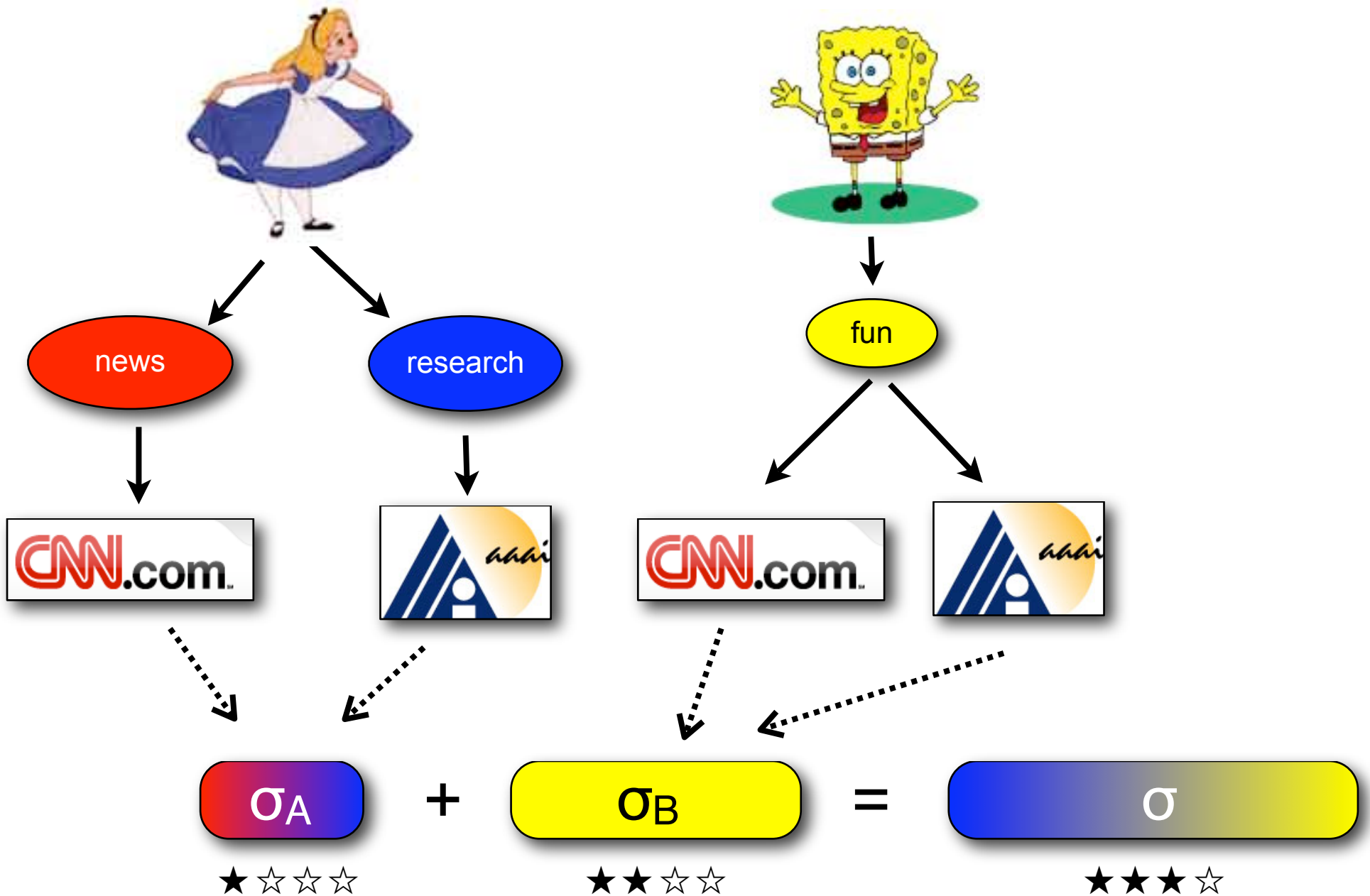
social





givealink.org

76788 links
1015 donations
117 registered users
Last updated: Apr 17, 2006





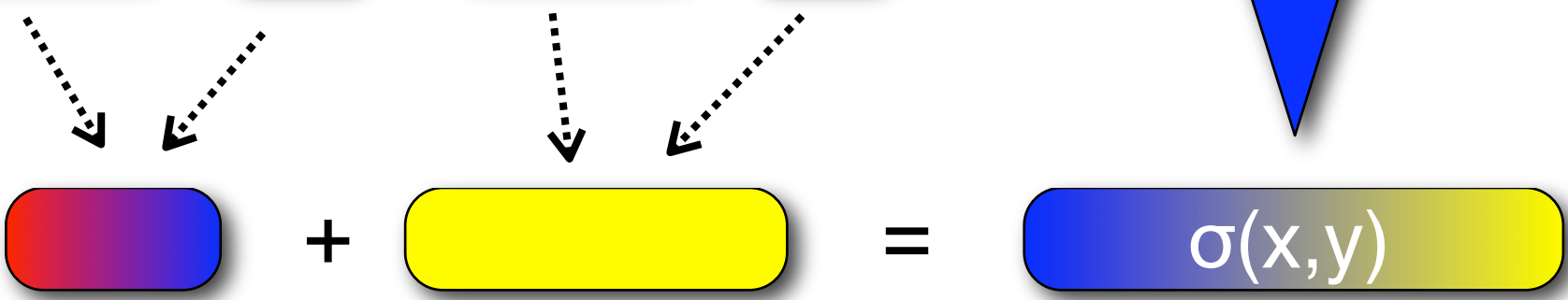
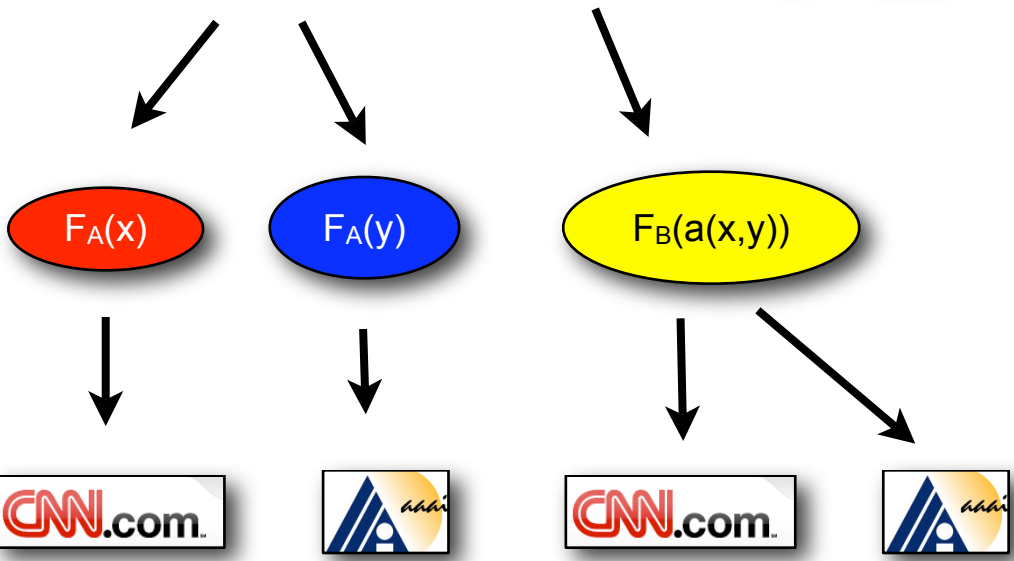
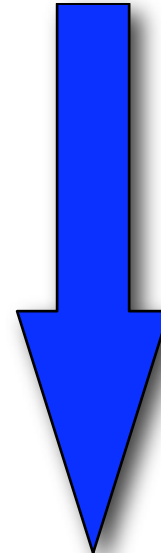
givealink.org

76788 links
1015 donations
117 registered users
Last updated: Apr 17, 2006



$\sigma(x, y)$

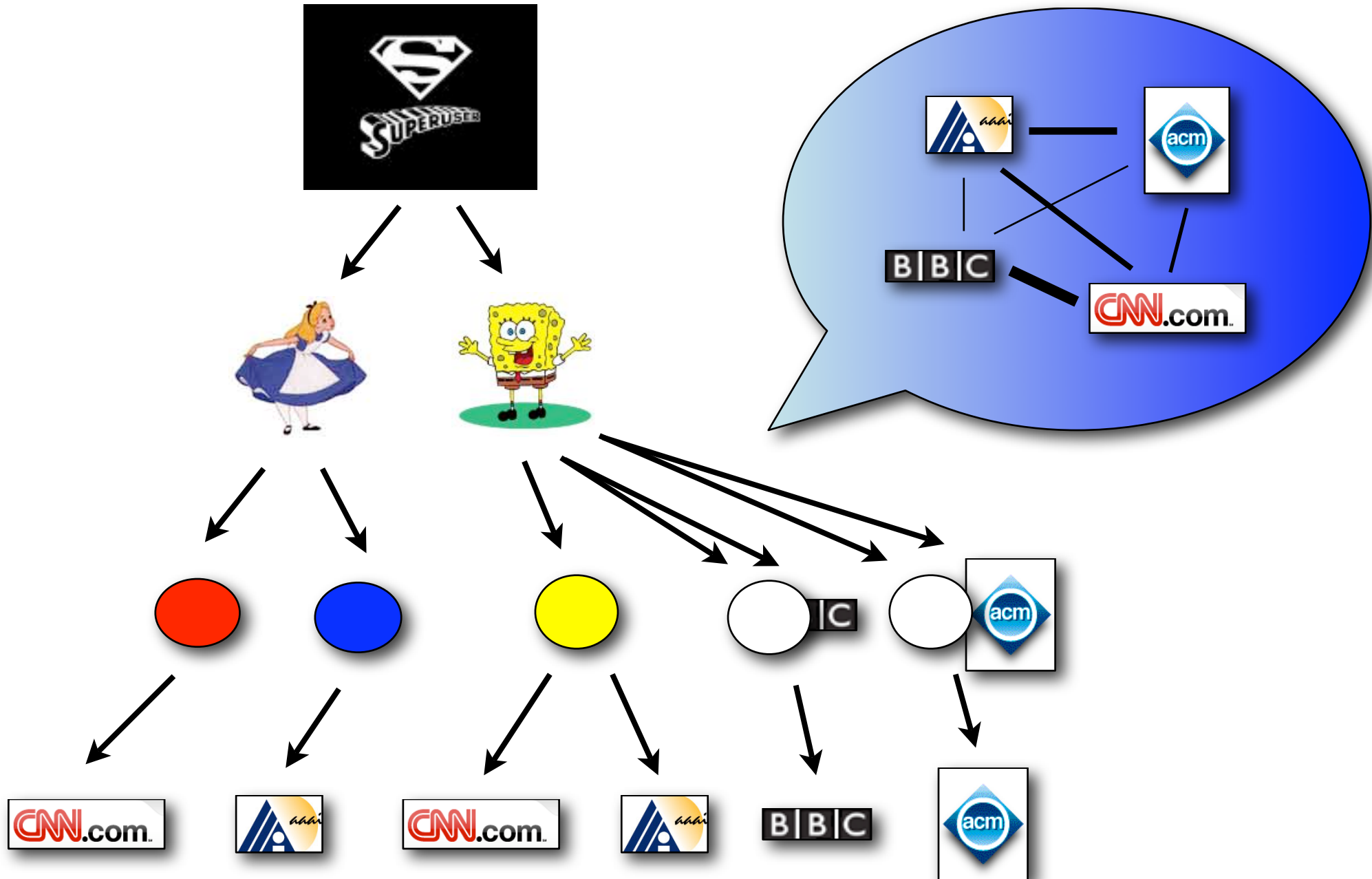
$$\sigma(x, y) = \frac{1}{N} \sum_{u=1}^N \frac{2 \log \left(\frac{|F_u[a(x, y)]|}{|R_u|} \right)}{\log \frac{|F_u(x)|}{|R_u|} + \log \frac{|F_u(y)|}{|R_u|}}$$





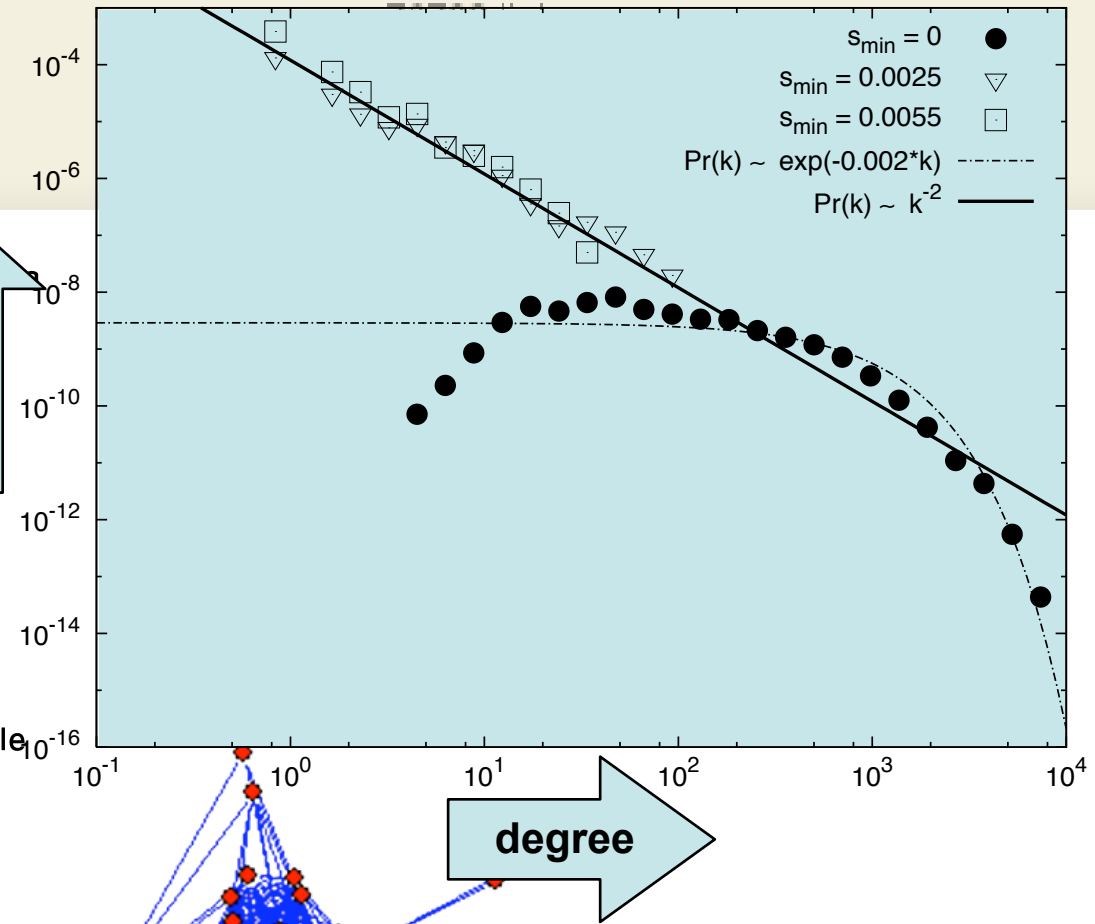
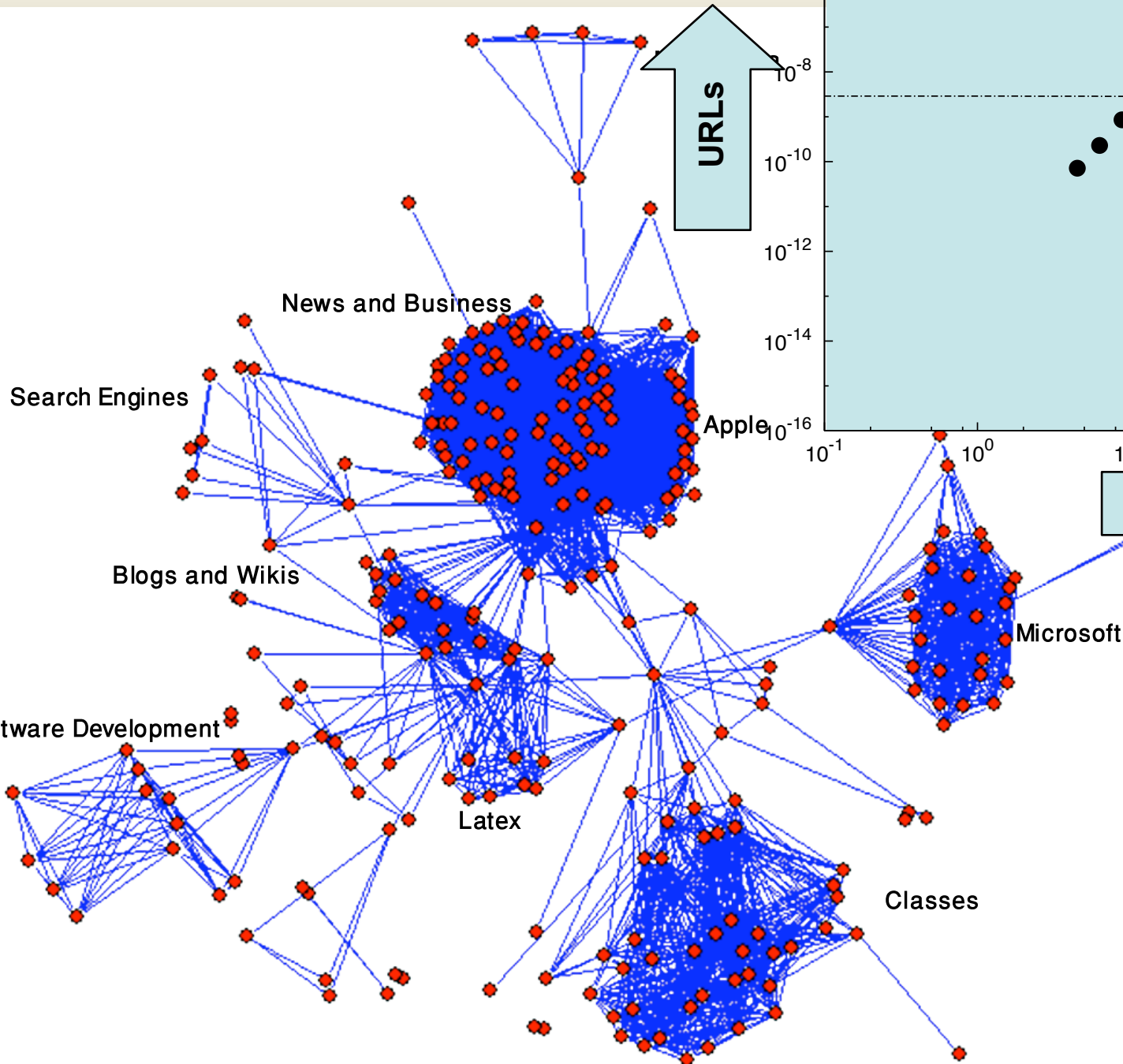
givealink.org

76788 links
1015 donations
117 registered users
Last updated: Apr 17, 2006





givealink.org





Link Recommender

Enter URL or keyword to get related bookmarks. ([More Options](#))

Sort by Bookmark: Similarity ? Novelty ? Prestige ?

Similar URLs sorted by Similarity

1-10 of 360 recommendations for <http://www.webir.org/>

<http://www.parc.xerox.com/istl/groups/lea/dynamics.shtml>

<http://www.parc.xerox.com/istl/groups/lea/dynamics.shtml>

http://www.springer.de/cgi-bin/search_book.pl?isbn=3-540-65112-8

http://www.springer.de/cgi-bin/search_book.pl?isbn=3-540-65112-8

[Web Research Collections - Web Track](#)

<http://es.csiro.au/TRECWeb/>

[DIMACS Workshop on Internet and WWW Measurement, Mapping and Modeling](#)

<http://dimacs.rutgers.edu/Workshops/Internet/>

<http://www.ibm.com/java/fetuccino/>

<http://www.ibm.com/java/fetuccino/>

[Web Term Document Frequency Form](#)

<http://elib.cs.berkeley.edu/docfreq/index.html>

[Finding Out About](#)

<http://www.cs.ucsd.edu/~rik/foa/>

[Bibliometrics of the World Wide Web: An Exploratory Analysis of the Intellectual Structure of Cyberspace](#)

<http://sherlock.berkeley.edu/asis96/asis96.html>

[Terabyte TREC Homepage](#)

<http://www-nlpir.nist.gov/projects/terabyte/>

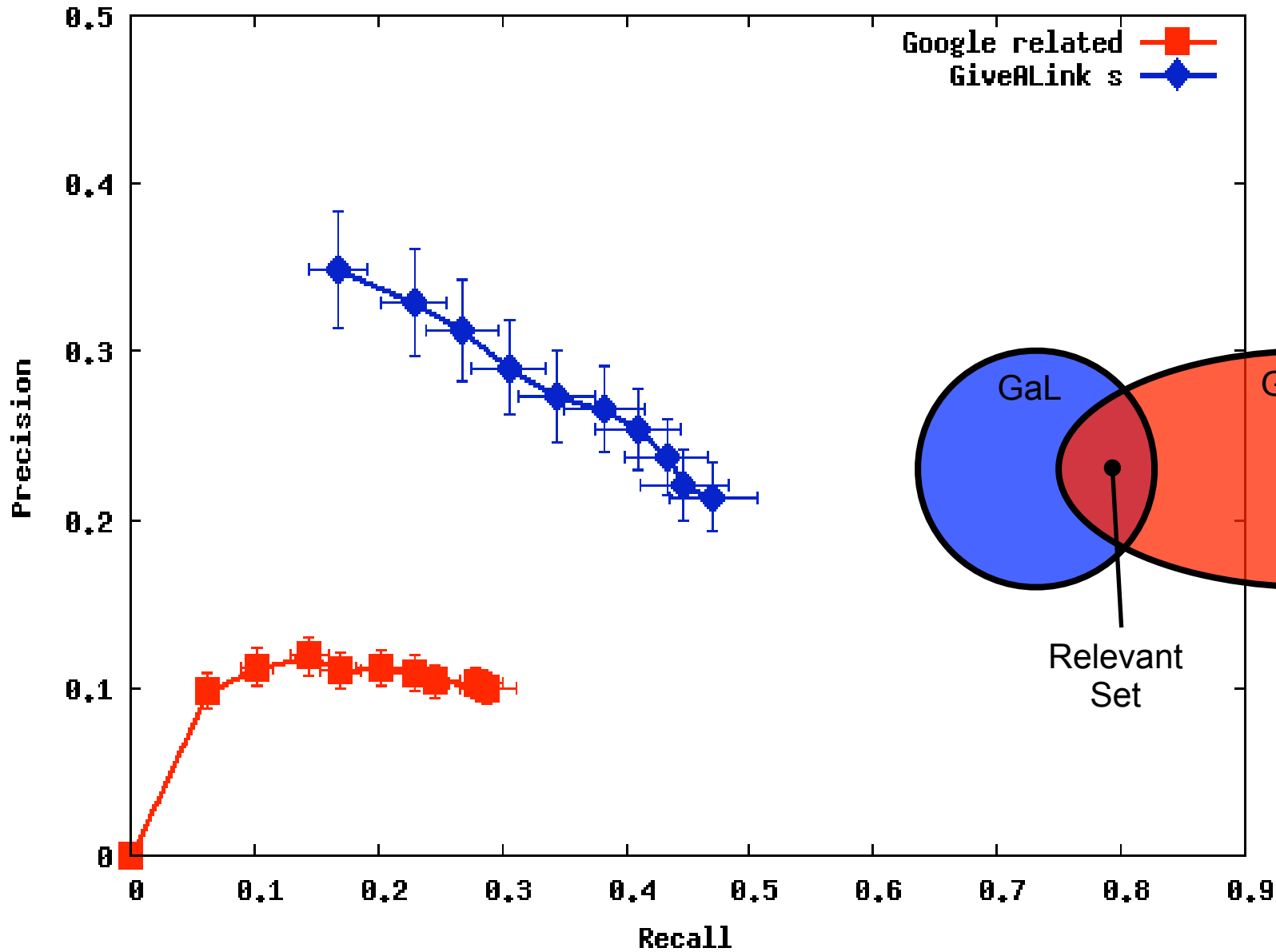
<http://www.neci.nj.nec.com/homepages/lawrence/websize.html>

<http://www.neci.nj.nec.com/homepages/lawrence/websize.html>

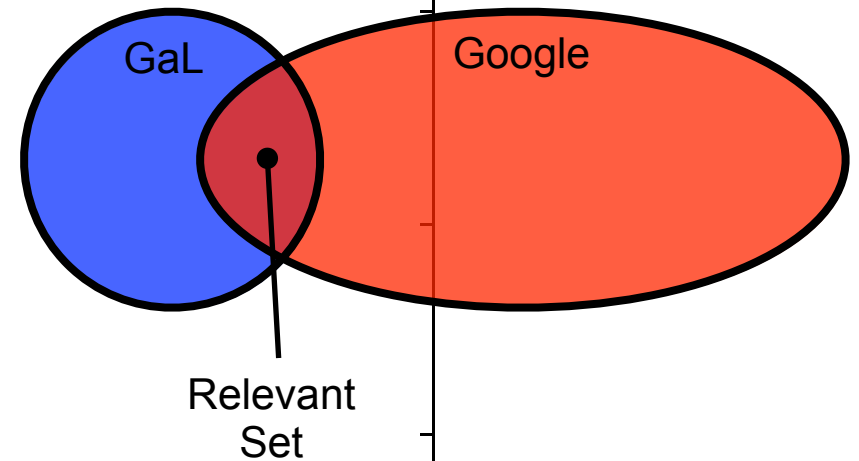


givealink.org

76788 links
1015 donations
117 registered users
Last updated: Apr 17, 2006



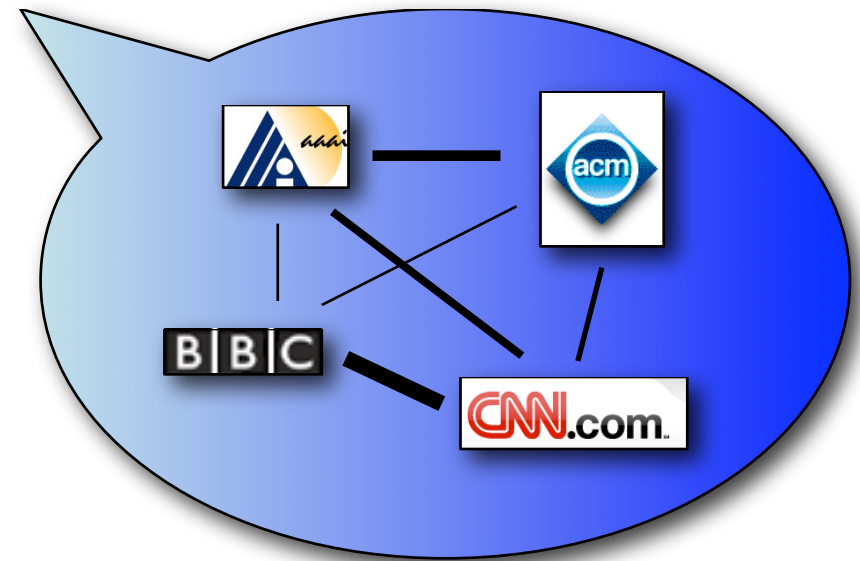
User Study





$$c(x) = \frac{1}{|U|} \sum_{y \in U} \left[1 + \min_{x \rightsquigarrow y} \sum_{(u,v) \in x \rightsquigarrow y} \left(\frac{1}{\sigma(u,v)} - 1 \right) \right]^{-1}$$

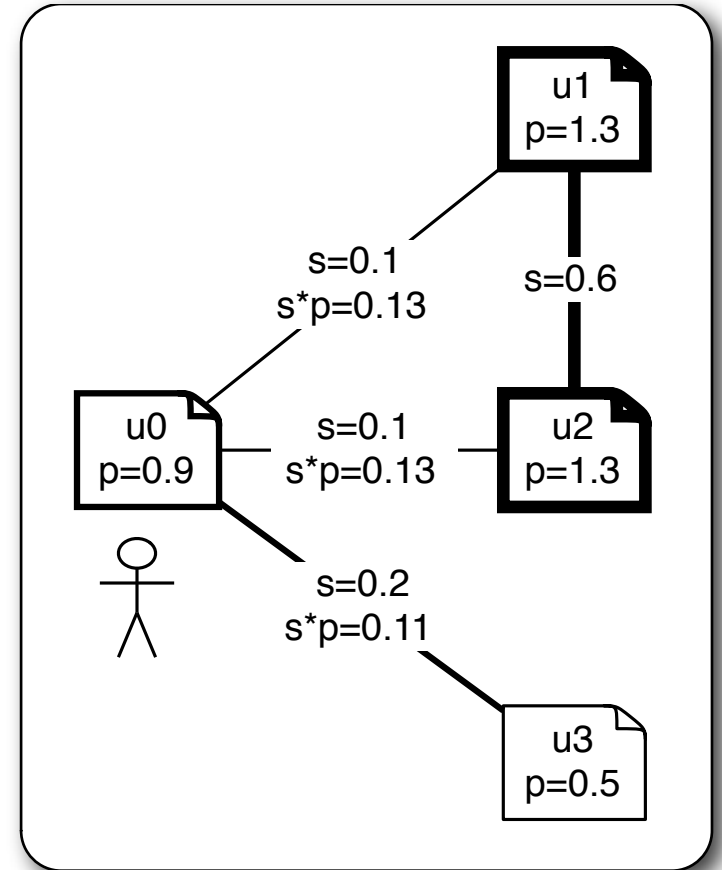
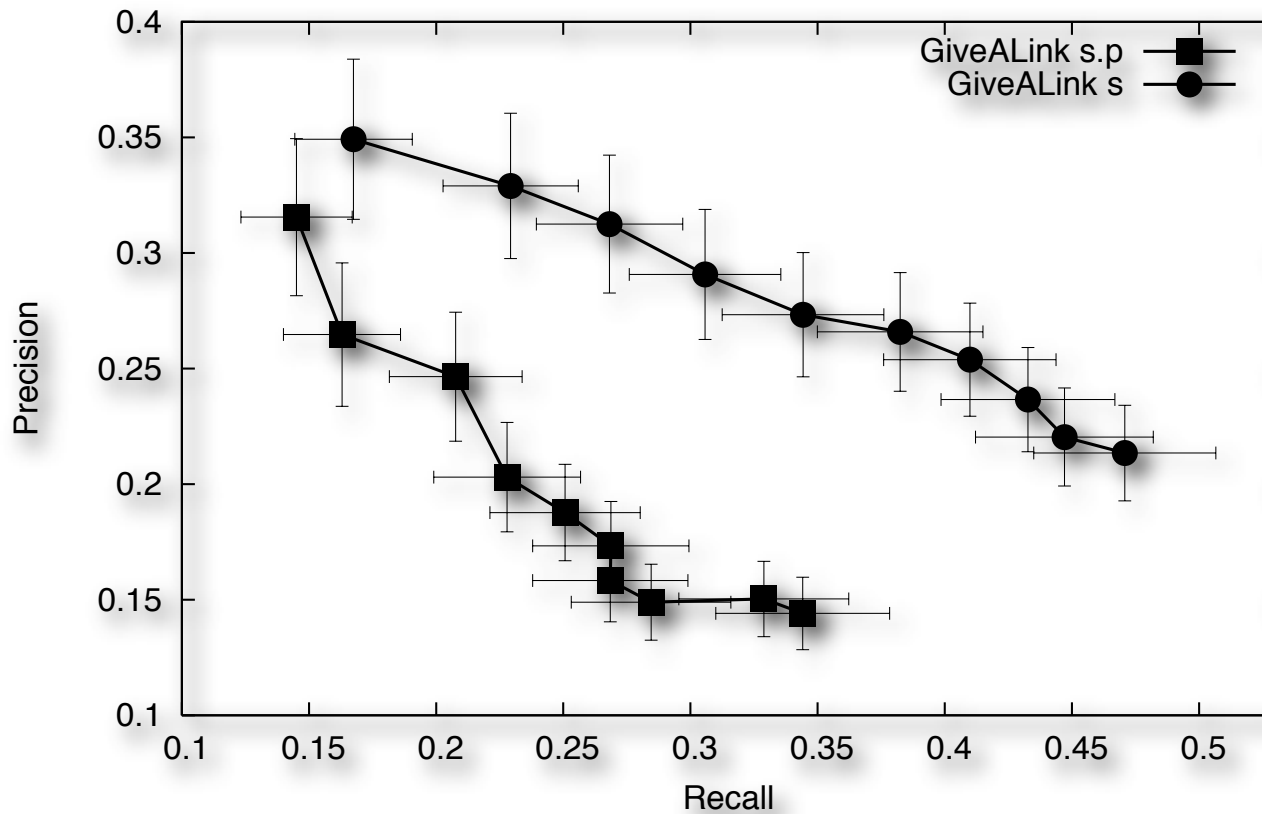
Centrality & Prestige



$$p_{t+1}(x) = (1 - \alpha) + \alpha \cdot \sum_{y \in U} \frac{\sigma(x, y) \cdot p_t(y)}{\sum_{z \in U} \sigma(y, z)}$$

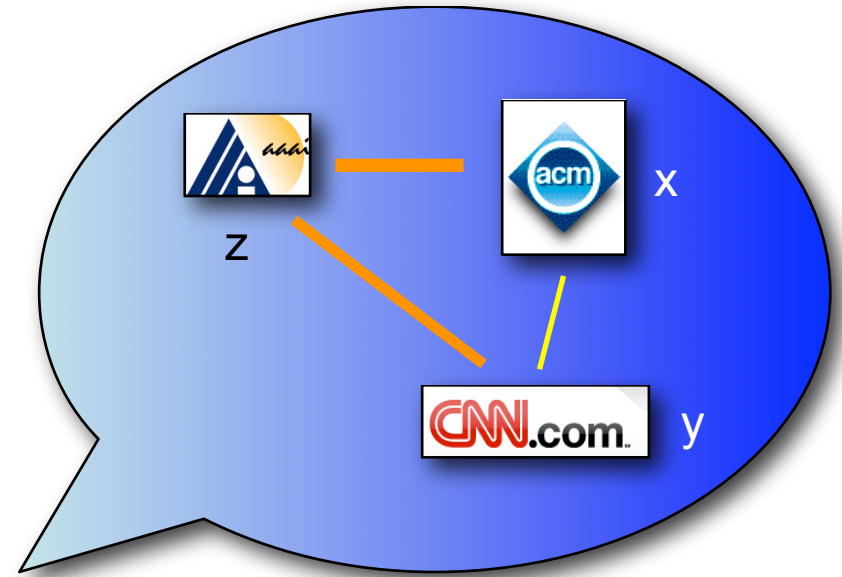


Ranking by Similarity * Prestige





Ranking & Recommendation by Novelty

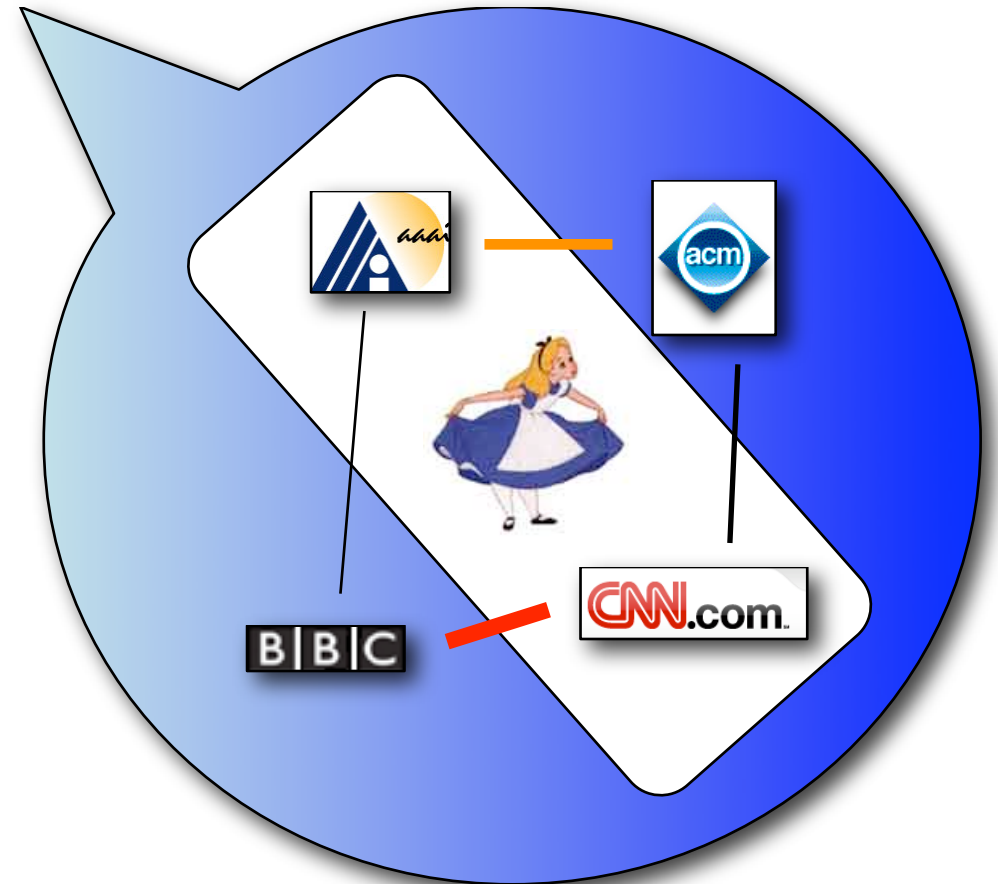


$$\nu(x, y) = \frac{\left[1 + \min_{z \in U} \left(\frac{1}{\sigma(x, z)} + \frac{1}{\sigma(z, y)} - 2 \right) \right]^{-1}}{\sigma(x, y)}$$



$$\eta(x, A) = \max_{y \in A} \left[\sigma(x, y) \cdot \log \left(\frac{N}{N(y)} \right) \right]$$

Personalization





collaborative filtering,
social semantic similarity,
unlinked pages,
multimedia content, trust

scalability,
density



spam

General

Surprise Me

Getting Started

1. [Register](#)
2. [Donate your bookmarks to science](#)
3. [Get personalized recommendations/search results](#)
Also [manage your bookmarks](#) or add a new URL while you browse
4. [Spread the word!](#)
More users means more data means better results!

Top Ten Bookmarks ?

1. <http://www.iuma.com/>
2. [Sonicnet.com](#)
3. [Error](#)
4. [the prog organ - progressive rock reviews](#)
5. [Hard Rock and Heavy Metal Radio - Video](#)
6. [GODS OF MUSIC - Music Reviews For The Independent Music Scene](#)
7. [RUTHLESS REVIEWS: MUSIC](#)
8. [Welcome! -- Rate Your Music](#)
9. [Tiny Mix Tapes](#)
10. [Reviews of Indie Albums on Irish music webzine CLUAS.com](#)

Stats

81949 links
3277880 relationships
1049 donations
184 registered users
Last Updated: Sat Aug 5 16:10:58 2006



Donate! givealink.org

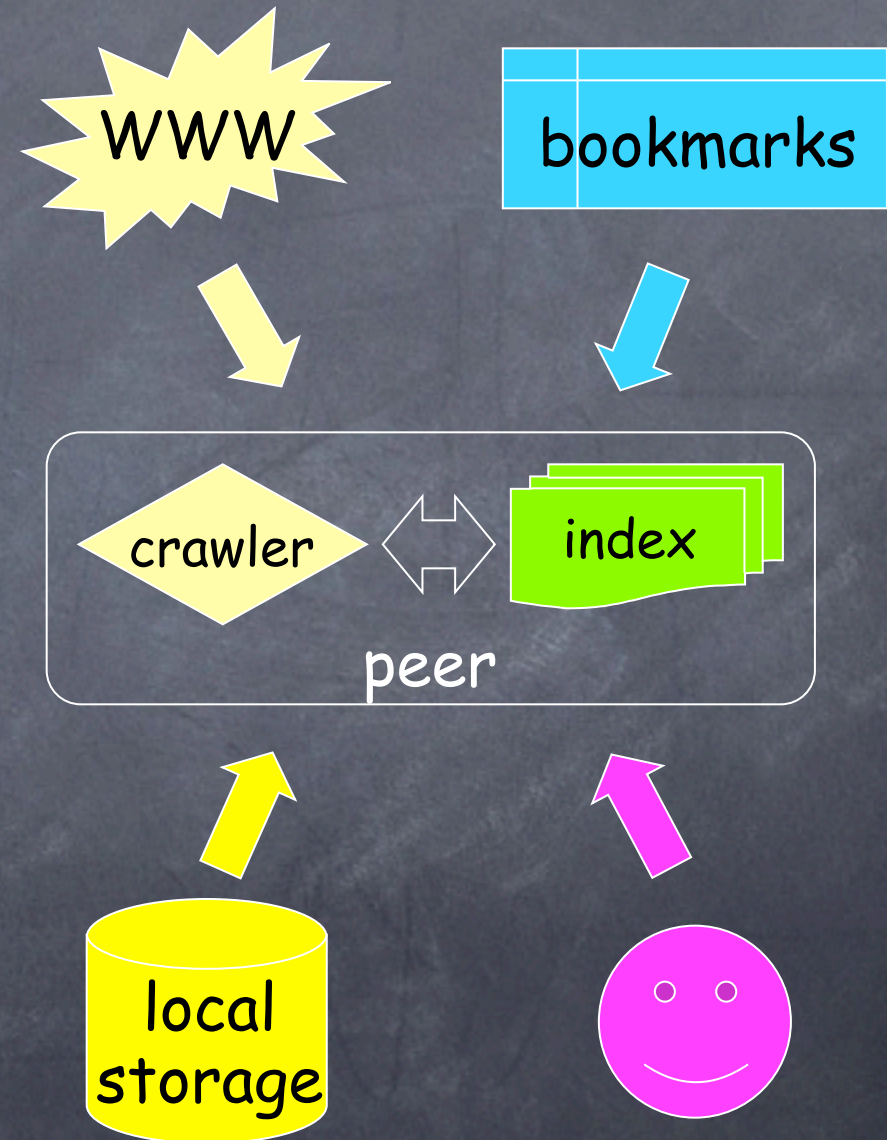


The GiveALink data is licensed under a [Creative Commons License](#).
© 2005, the Trustees of [Indiana University](#)

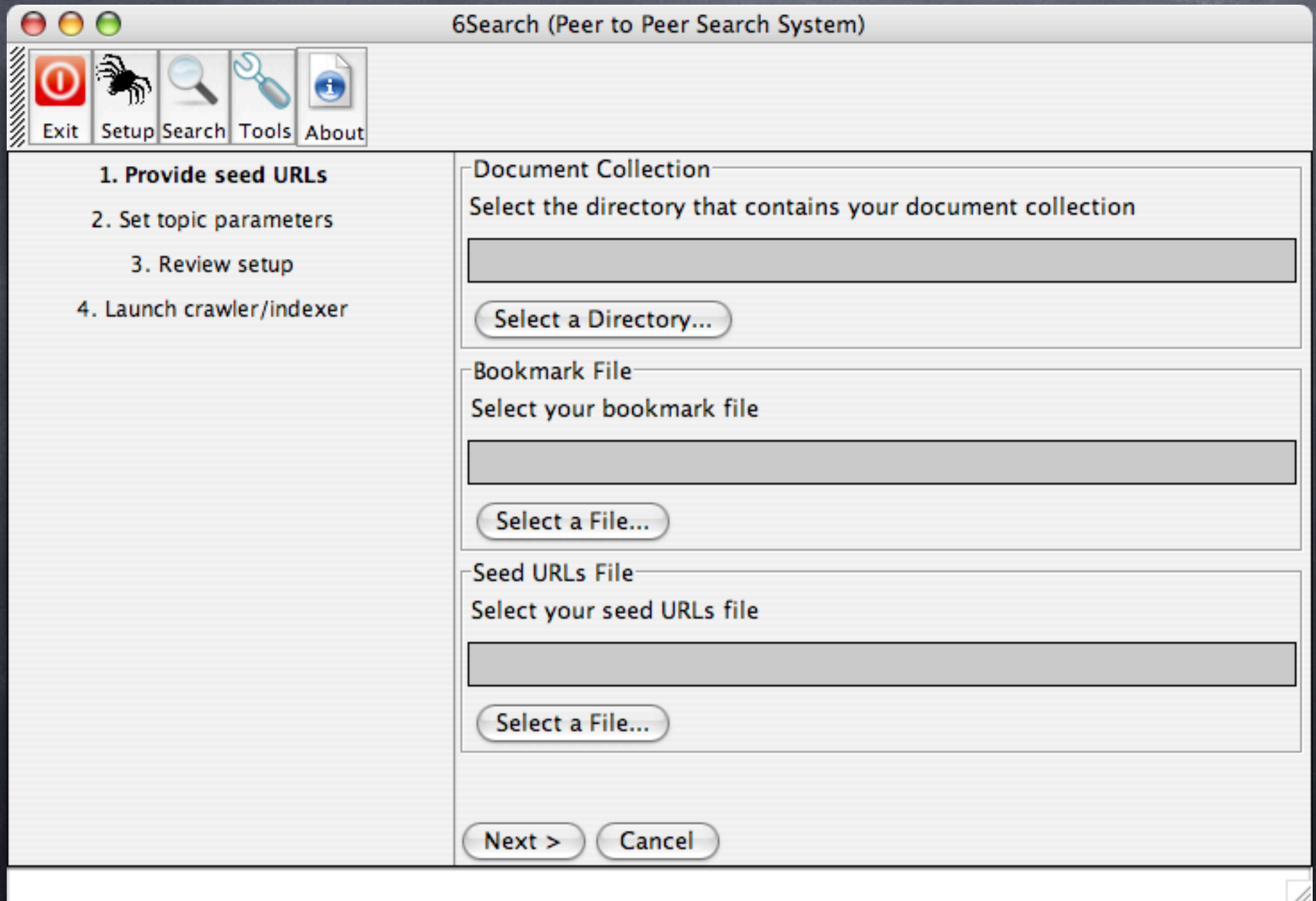
Discussion: Critical mass?

Peer distributed crawling and collaborative searching

- **Crawl** when idle
 - Start from **bookmarks**
 - **Past queries or selected documents** as topics
- **Index** locally
 - User **relevance feedback** is natural
 - No centralized coordination bottleneck
 - Unlike grub.org, hyperbee.com



<http://homer.informatics.indiana.edu/~nan/6S/>



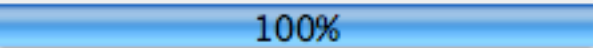
<http://homer.informatics.indiana.edu/~nan/6S/>

6Search (Peer to Peer Search System)

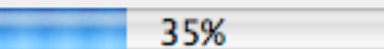
Exit Setup Search Tools About

1. Provide seed URLs
2. Set topic parameters
3. Review setup
- 4. Launch crawler/indexer**

Crawling Progress

11/08 17:30:01  100%

Indexing Progress

11/08 17:31:06  35%

Information Summary

Please verify the information you entered.

Document Collection:
/Users/fil/Documents/Homepages/Informatics

Bookmark File:
/Users/fil/Library/Safari/Bookmarks.plist

Crawling Topic:
web mining search peer crawl networks agents modeling

Seed URLs File:

Number of Pages:
100

Your Email:
fil@indiana.edu

Top for URLs:

<http://homer.informatics.indiana.edu/~nan/6S/>

6Search (Peer to Peer Search System)

Exit Setup Search Tools About

Indexing Date	Document Title	Document Url	Document Tag
2006/11/08 ...	SiliconBeat: News about tech money and in...	http://www.siliconbeat.com/	web mining search pe...
2006/11/08 ...	Investment Management Research Services	http://www.round.table.com/i...	web mining search pe...
2006/11/08 ...	VentureBeat » Bold start-up, Powerset, ab...	http://venturebeat.com/2006/...	web mining search pe...
2006/11/08 ...	Yahoo!	http://www.yahoo.com/	web mining search pe...
2006/11/08 ...	The J Curve	http://jurvetson.blogspot.com/	web mining search pe...
2006/11/08 ...		http://www.webconferences.co...	web mining search pe...
2006/11/08 ...	Search Boards for threads - Forum search ...	http://www.boardtracker.com/...	web mining search pe...
2006/11/08 ...	Techcrunch » Blog Archive » Will Powerset ...	http://www.techcrunch.com/20...	web mining search pe...
2006/11/08 ...	Jeff Clavier	http://softtechvc.blogs.com/ab...	web mining search pe...
2006/11/08 ...	Challengelist	http://www.research.att.com/%...	web mining search pe...
2006/11/08 ...		http://www.gutenberg.org/dirs...	web mining search pe...
2006/11/08 ...		http://www.w3.org/WAI/UA/W...	web mining search pe...
2006/11/08 ...	My AOL	http://feeds.my.aol.com/?url=...	web mining search pe...
2006/11/08 ...	Donor Suitability Workshop: Donor History ...	http://www.fda.gov/cber/minu...	web mining search pe...
2006/11/08 ...	ODP - Open Directory Project	http://www.dmoz.org/	web mining search pe...
2006/11/08 ...	The CC Chemokine Thymus-derived Chem...	http://www.jem.org/cgi/conten...	web mining search pe...
2006/11/08 ...	absintelagent : Messages : 258-287 of 287	http://tech.groups.yahoo.com/...	web mining search pe...
2006/11/08 ...	Survey: Information Gathering and Knowled...	http://www.cio.com/research/...	web mining search pe...
2006/11/08 ...	British Blogs	http://www.britishblogs.co.uk/	web mining search pe...
2006/11/08 ...	NJDEP Compliance & Enforcement - A decl...	http://www.state.nj.us/dep/en...	web mining search pe...
2006/11/08 ...	Barney Pell's Weblog: Search Archives	http://www.barneypell.com/ar...	web mining search pe...
2006/11/08 ...	Flatable Bradley Horowitz	http://www.elatable.com/blog/	web mining search pe...

Advanced Tool

Delete row Undelete All Delete All

<http://homer.informatics.indiana.edu/~nan/6S/>

6Search (Peer to Peer Search System)

Exit Setup Search Tools About

6search web mining Search

Local Search Only Yahoo Search Google Search Peer network connected

6S search Google search Yahoo search

[Data Mining - Home Page \(Misc\)](#)
... Data Mining Software (index) Data Mining Events (index) Data Mining General/Misc (index) People working ... OnLine Analytical Processing (OLAP) , Data ...
<http://www.the-data-mine.com/> [Similar pages \(Power by GiveALink\)](#)
Contributors: wls_iceman#1,wls_angel#3,wls_beast#7

[KDnuggets: Data Mining, Web Mining, and Knowledge Discovery Guide](#)
... Data Mining, Web Mining, Text Mining, and Knowledge Discovery ... KDnuggets: Data Mining, Web Mining, and Knowledge Discovery ...
<http://www.kdnuggets.com/> [Similar pages \(Power by GiveALink\)](#)
Contributors: wls_iceman#1,wls_angel#3

[Data Mining Software in the Yahoo! Directory](#)
... Data Mining Software in the ...
http://dir.yahoo.com/Business_and_Economy/Business_to_Business/Computers/Software/Databases/Data_Mining/ [Similar pages \(Power by GiveALink\)](#)
Contributors: wls_iceman#1,wls_beast#7

[UMBC Agent Web -- news and information on software agent technology](#)
... UMBC Agent Web -- news and information ...
http://agents.umbc.edu/Topics/Related_Topics/Information_retrieval_and_knowledge_management/ [Similar pages \(Power by GiveALink\)](#)
Contributors: wls_iceman#1,wls_beast#7

Done

<http://homer.informatics.indiana.edu/~nan/6S/>

The screenshot shows a web browser window titled "The search results of 6S system". The address bar displays "http://localhost:1999/result/" and the search bar contains "6S web mining". A dropdown menu is open, listing search engines: Google, Yahoo, 6S 6Search (local), 6S 6Search (peers) (highlighted), Amazon.com, Answers.com, Creative Commons, eBay, and Manage Search Engines... The main content area lists search results for "Data Mining, Web Mining, and Knowledge Discovery Guide", "Data Mining - Home Page (Misc)", "Data Mining Software in the Yahoo! Directory", "UMBC Agent Web -- news and information on software agent technology", and "BUBL LINK: Information retrieval".

KDnuggets: Data Mining, Web Mining, and Knowledge Discovery Guide
Newsletter on the data mining and knowledge industries, offering information on data mining, know
and web mining software, courses, jobs, publications, and meetings.
<http://www.kdnuggets.com/>
Contributors: Main6S#0,wls_angel#3,wls_iceman#1
[Similar pages \(Power by GiveALink\)](#)

Data Mining - Home Page (Misc)
... in an unrestricted hands-on **web**. Software Information on **Data Mining** Software Evaluate
Analytical Processing (OLAP) , **Data Mining** ...
<http://www.the-data-mine.com/>
Contributors: Unknown#7,wls_angel#3,wls_iceman#1
[Similar pages \(Power by GiveALink\)](#)

Data Mining Software in the Yahoo! Directory
... for IT professionals focusing on **data mining**, **data** analysis, and reports ... to Business > Computers > Software > Databases >
Data ...
http://dir.yahoo.com/Business_and_Economy/Business_to_Business/Computers/Software/Databases/Data_Mining/
Contributors: Unknown#7,wls_iceman#1
[Similar pages \(Power by GiveALink\)](#)

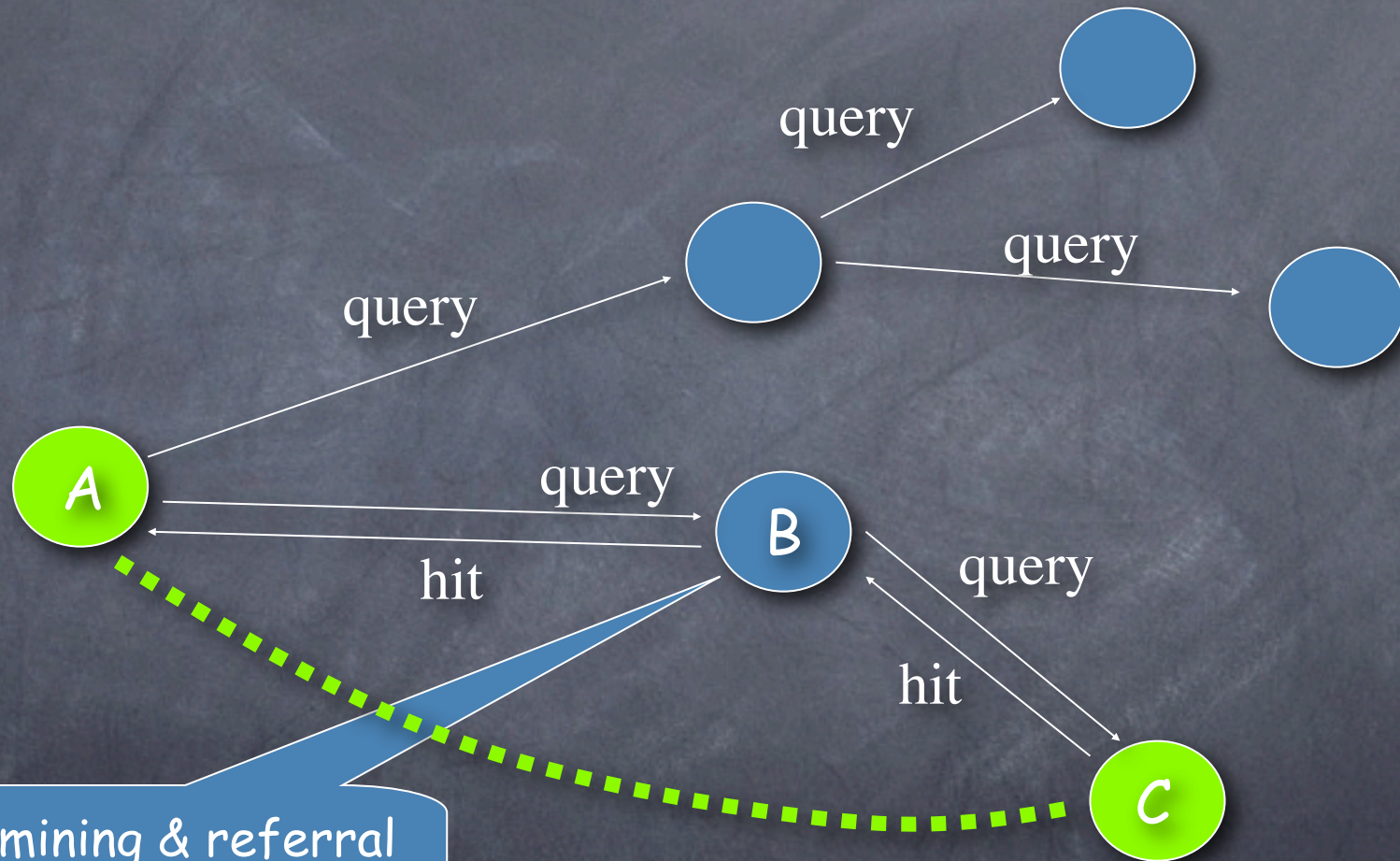
UMBC Agent Web -- news and information on software agent technology
... UMBC Agent **Web** -- news and information ...
http://agents.umbc.edu/Topics/Related_Topics/Information_retrieval_and_knowledge_management/
Contributors: Unknown#7,wls_iceman#1
[Similar pages \(Power by GiveALink\)](#)

BUBL LINK: Information retrieval
... and Forward Knowledge Approach **Data Mine: Data Mining** and Knowledge Discovery ... the Gaps KD Mine: **Data** ...
<http://bubl.ac.uk/link/i/informationretrieval.htm>
Contributors: Unknown#7,wls_iceman#1
[Similar pages \(Power by GiveALink\)](#)

Library and Information Science > Information Retrieval in the Yahoo! Directory

Done

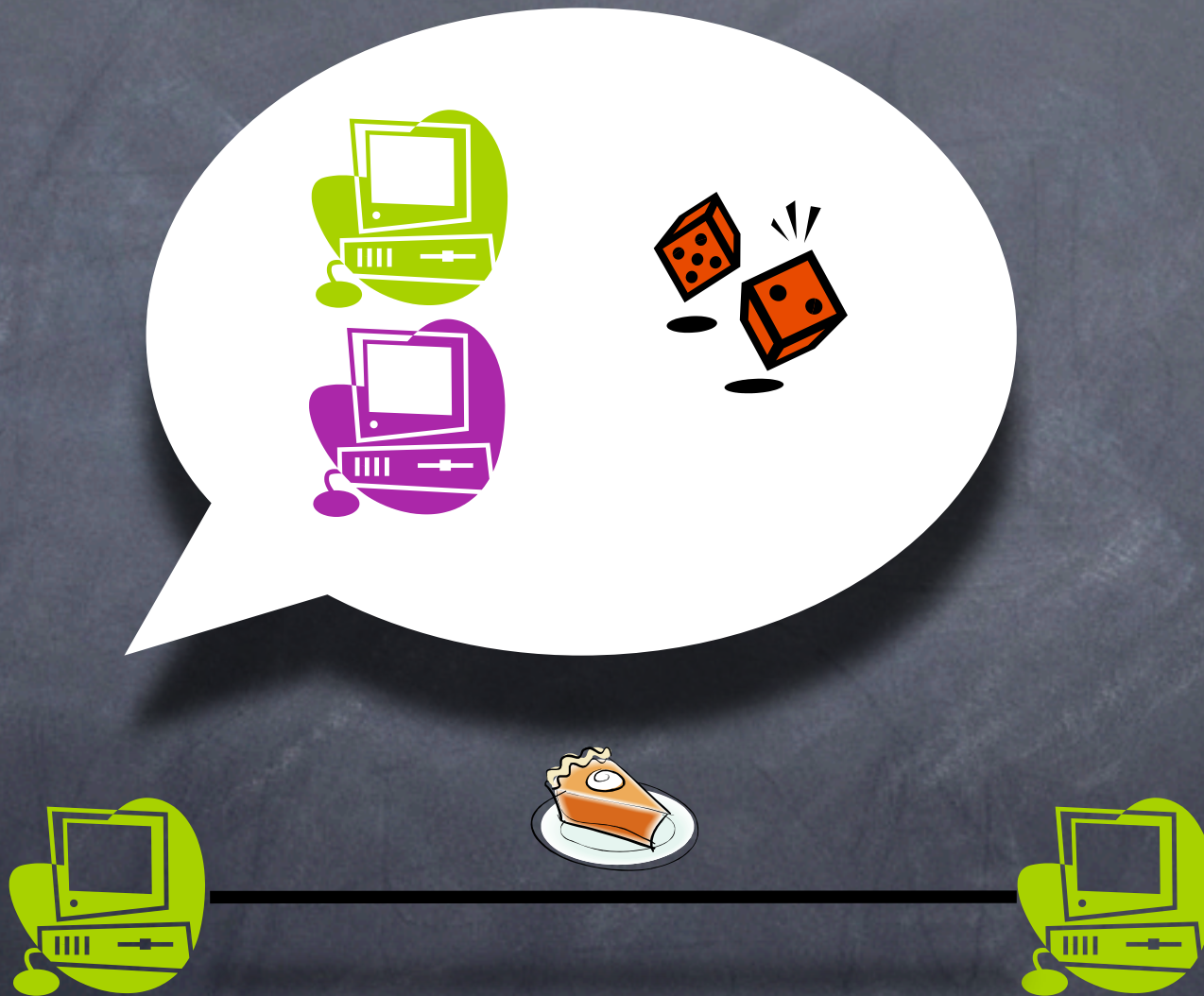
6S: peer distributed crawling and collaborative searching



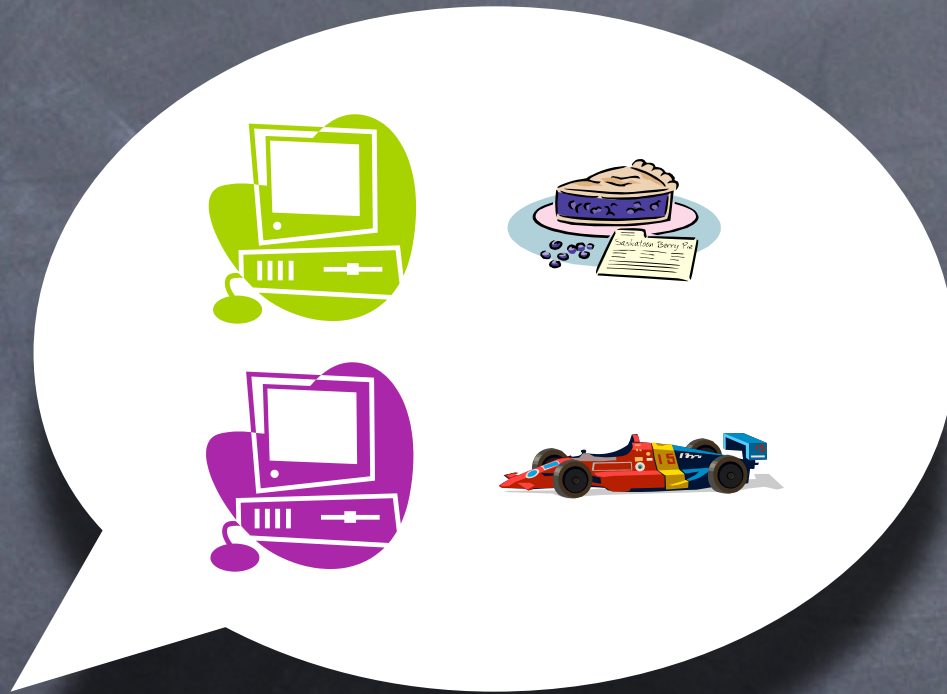
Data mining & referral opportunities

Emerging communities

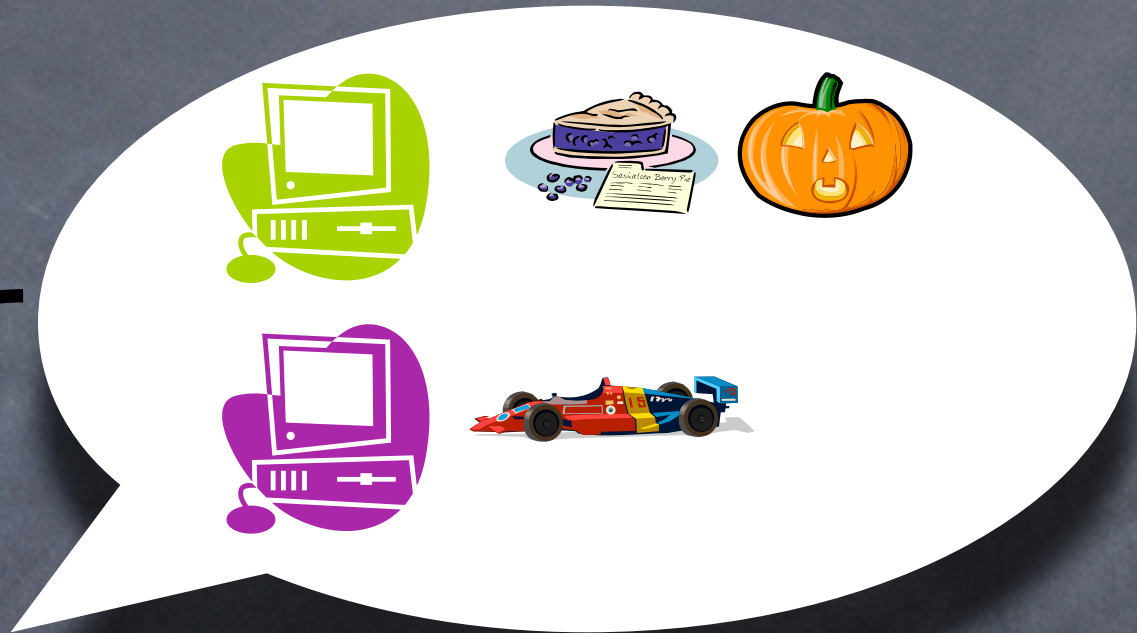
Algorithm 1: Random Known



Algorithm 2: Greedy



Algorithm 3: Reinforcement



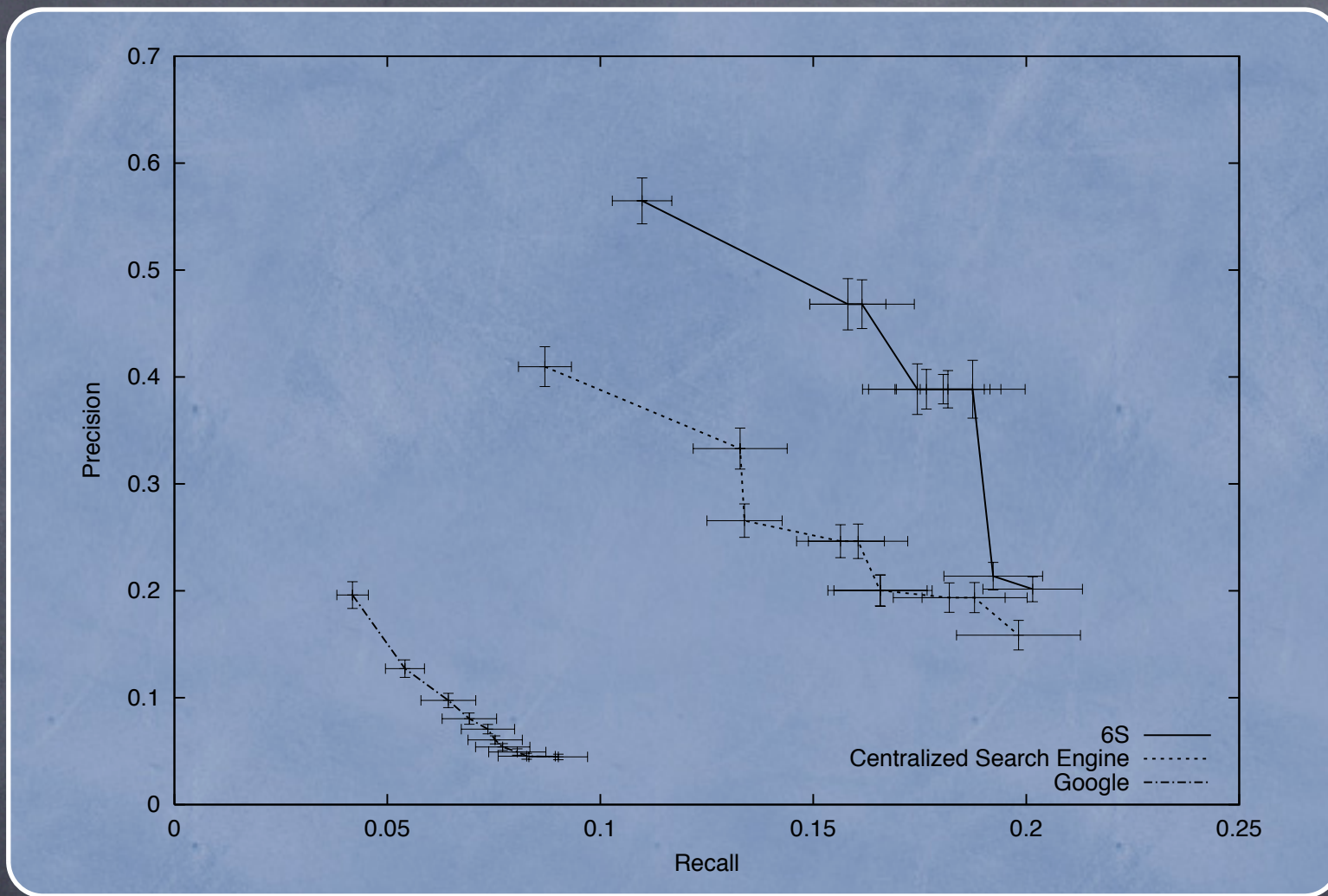
Learning Peers

- A focused profile W^f of other peers for terms in queries
- An expanded profile W^e of other peers for terms not in queries, but that co-occur with query terms and with higher frequency

$$w_{i,p}(t+1) = (1 - \gamma) \cdot w_{i,p}(t) + \gamma \cdot \left(\frac{S_p + 1}{S_l + 1} - 1 \right)$$

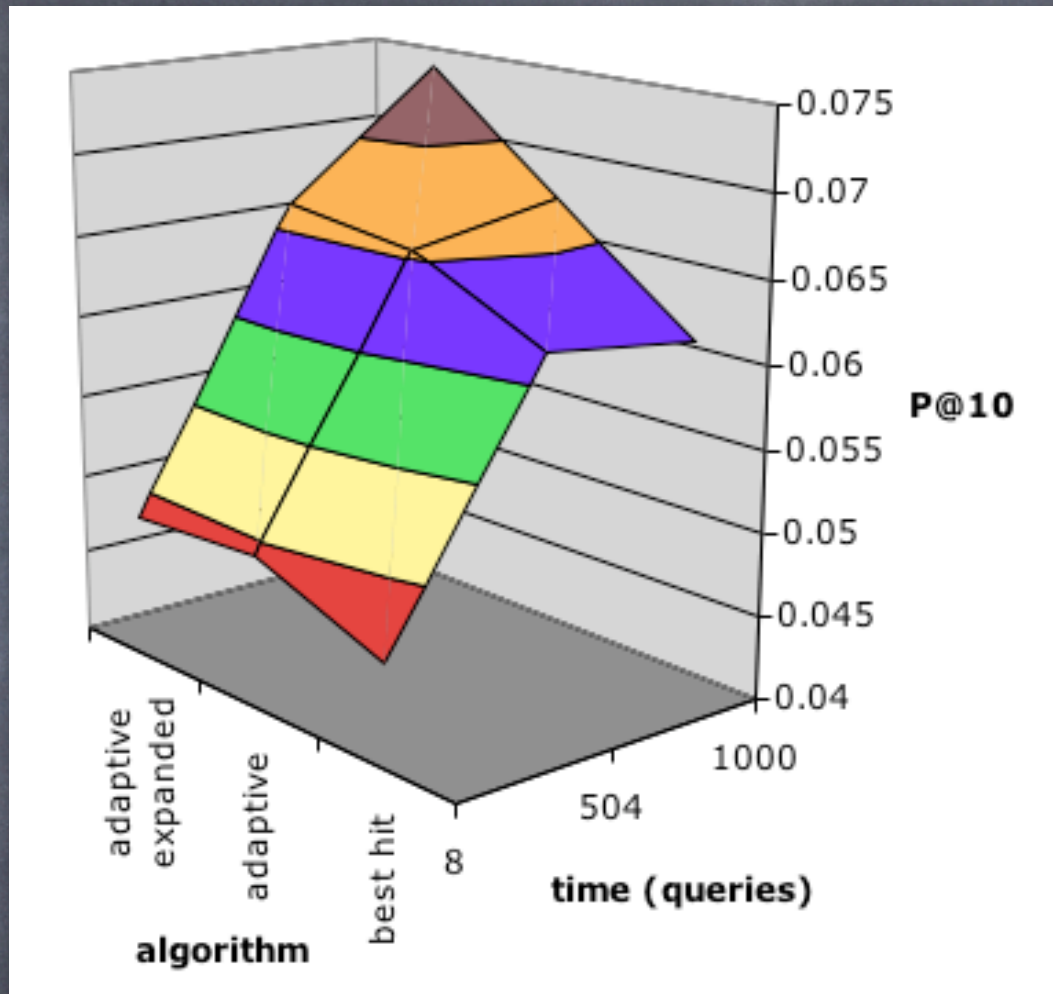
- Query routing:

$$\sigma(p, Q) = \sum_{i \in Q} \left[\alpha \cdot w_{i,p}^f + (1 - \alpha) \cdot w_{i,p}^e \right]$$

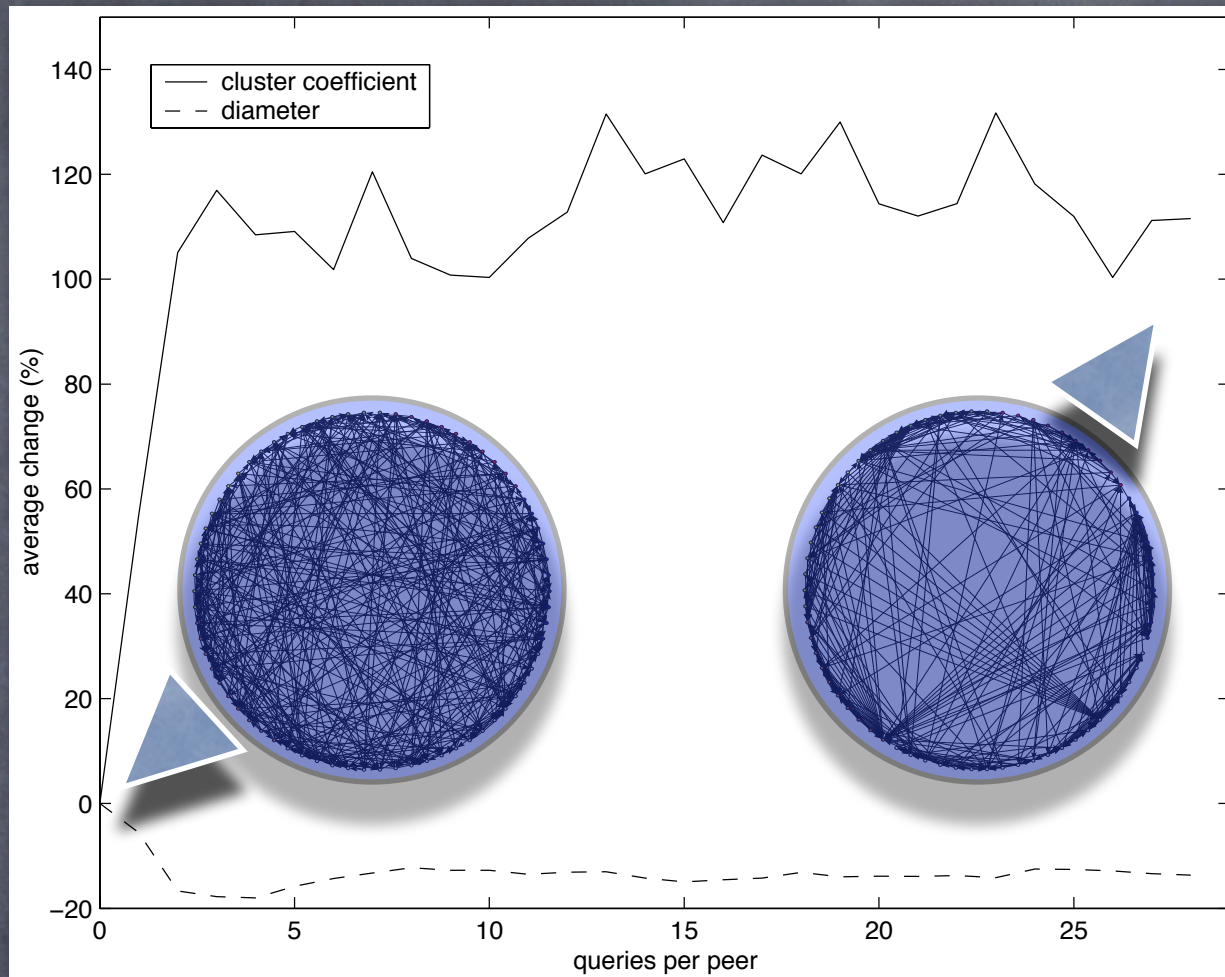


Distributed vs Centralized

WTAS 2005

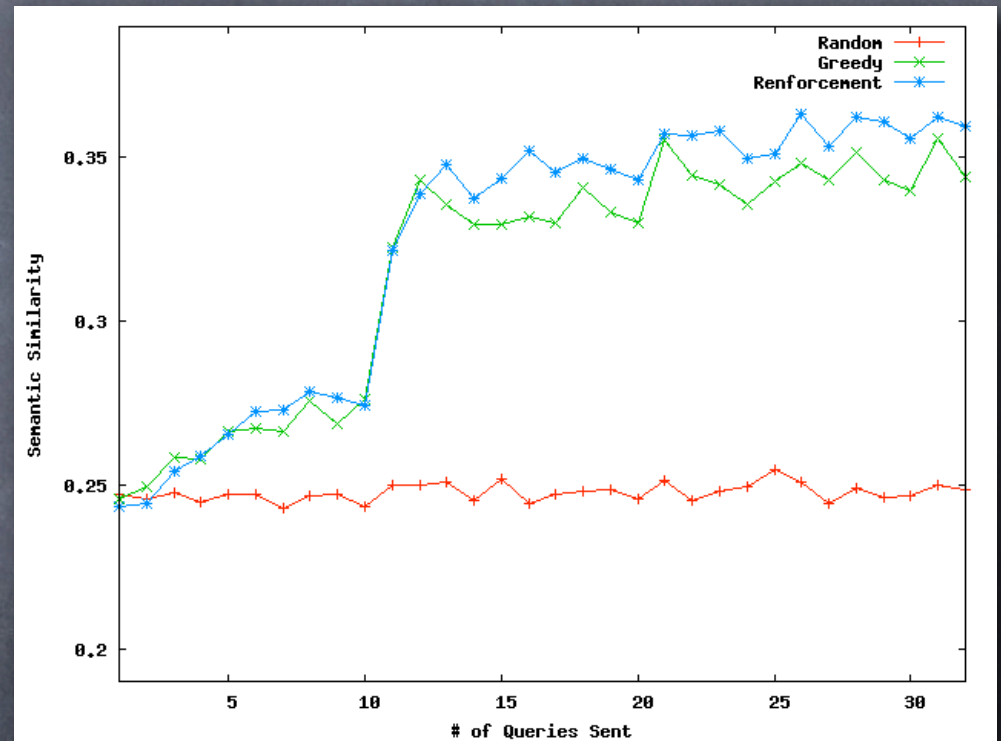


P@10



Small-world

Semantic similarity



Discussion:

Would you use 6S?

Why or why not?

Thank you! Questions?

informatics.indiana.edu/fil

homer.informatics.indiana.edu/~nan



Research supported by NSF
CAREER Award IIS-0348940



08203848509823409582093845082830495
08 8666 734 10598 082 8509
0E 8485C 14 7093b450E 0495
3E 92830 106 79872E 2897
0E 86658 14234 782C 8509
0E 8485C 14 85820c 150E 0495
3E 9283 1680 123234 872E 2897
0E 76234 1680 123234 872E 2897

Data and Search Institute

Next generation tools and ideas for searching and manipulating large volumes of data

www.dataandsearch.org