

<http://poloclub.gatech.edu/cse6242>

CSE6242 / CX4242: **Data** & **Visual** Analytics

# Graphs / Networks

Centrality measures, algorithms, interactive applications

Duen Horng (Polo) Chau

Assistant Professor

Associate Director, MS Analytics

Georgia Tech

Partly based on materials by

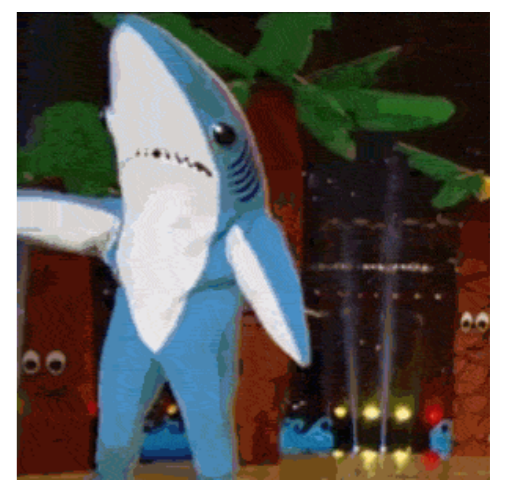
Professors Guy Lebanon, Jeffrey Heer, John Stasko, Christos Faloutsos, Parishit Ram (GT PhD alum; SkyTree), Alex Gray

**Centrality**  
= “Importance”

# Why Node Centrality?

What can we do if we can rank all the nodes in a graph (e.g., Facebook, LinkedIn, Twitter)?

- Find **celebrities** or influential people in a social network (Twitter)
- Find “**gatekeepers**” who connect communities (headhunters love to find them on LinkedIn)
- What else?



# More generally

Helps **graph analysis, visualization, understanding**, e.g.,

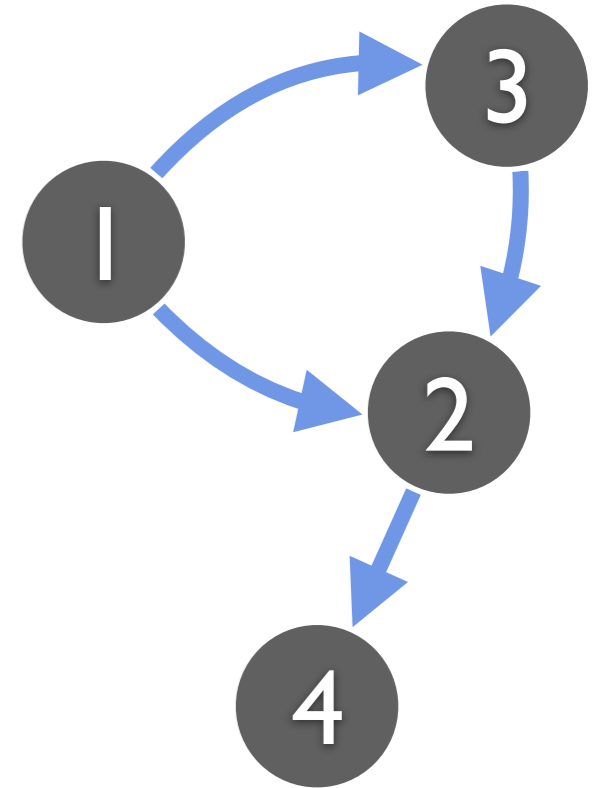
- Let us **rank** nodes, group or study them by centrality
- Only show subgraph formed by the **top 100 nodes**, out of the millions in the full graph
- **Similar to google search results** (ranked, and they only show you 10 per page)
- Most graph analysis packages already have centrality algorithms implemented. **Use them!**

Can also compute edge centrality.  
Here we focus on node centrality.

# Degree Centrality (easiest)

**Degree = number of neighbors**

- For directed graphs
  - **In degree** = No. of incoming edges
  - **Out degree** = No. of outgoing edges
- For undirected graphs, **only degree is defined.**
- Algorithms?
  - Sequential scan through **edge list**
  - What about for a **graph stored in SQLite?**



# Computing Degrees using SQL

Recall simplest way to store a graph in SQLite:

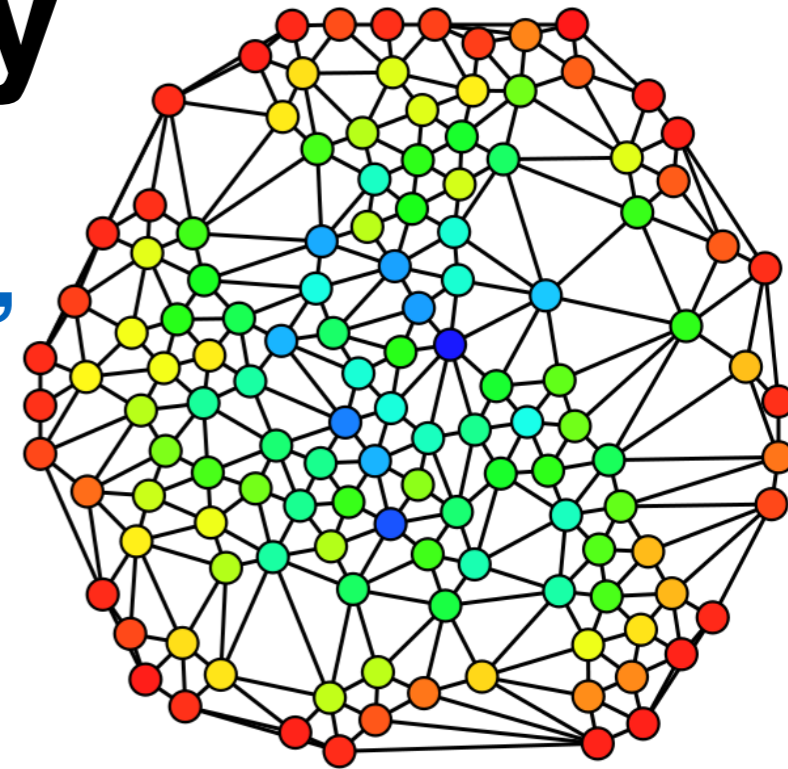
```
edges(source_id, target_id)
```

1. If slow, first create index for each column
2. Use **group by** statement to find **in degrees**

```
select count(*) from edges group by source_id;
```

# Betweenness Centrality

High betweenness = “gatekeeper”



Betweenness of a node  $v$

$$= \sum_{s \neq v \neq t \in V} \frac{\sigma_{st}(v)}{\sigma_{st}}$$

Number of shortest paths between  $s$  and  $t$  that **goes through  $v$**

Number of shortest paths between  $s$  and  $t$

= how often a node serves as the “bridge” that connects two other nodes.

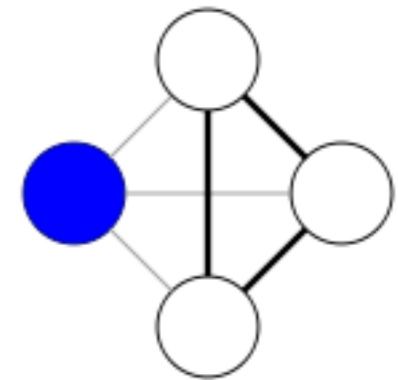
# (Local) Clustering Coefficient

A node's clustering coefficient is a measure of **how close the node's neighbors are from forming a clique.**

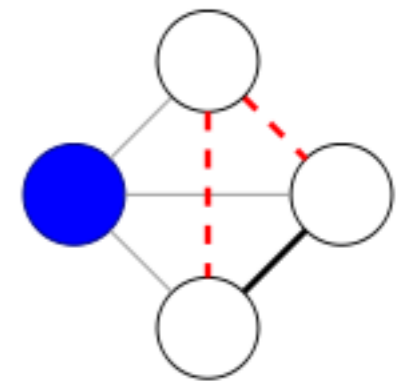
- 1 = neighbors form a clique
- 0 = No edges among neighbors

(Assuming undirected graph)

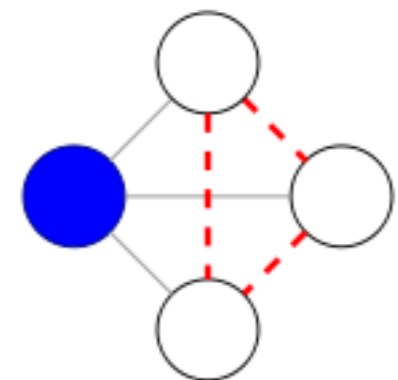
“Local” means it's for a node; can also compute a graph's “global” coefficient



$$c = 1$$



$$c = 1/3$$



$$c = 0$$



# Computing Clustering Coefficients...

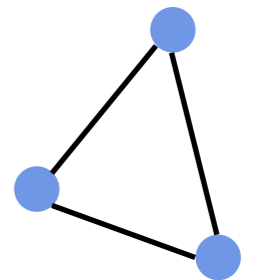
Requires **triangle counting**

Real social networks have a lot of triangles

- Friends of friends are friends

Triangles are **expensive** to compute

(neighborhood intersections; several approx. algos)



## Can we do that quickly?

Algorithm details:

Faster Clustering Coefficient Using Vertex Covers

<http://www.cc.gatech.edu/~ogreen3/docs/2013VertexCoverClusteringCoefficients.pdf>

# Super Fast Triangle Counting

## [Tsourakakis ICDM 2008]



But: triangles are expensive to compute  
(3-way join; several approx. algos)

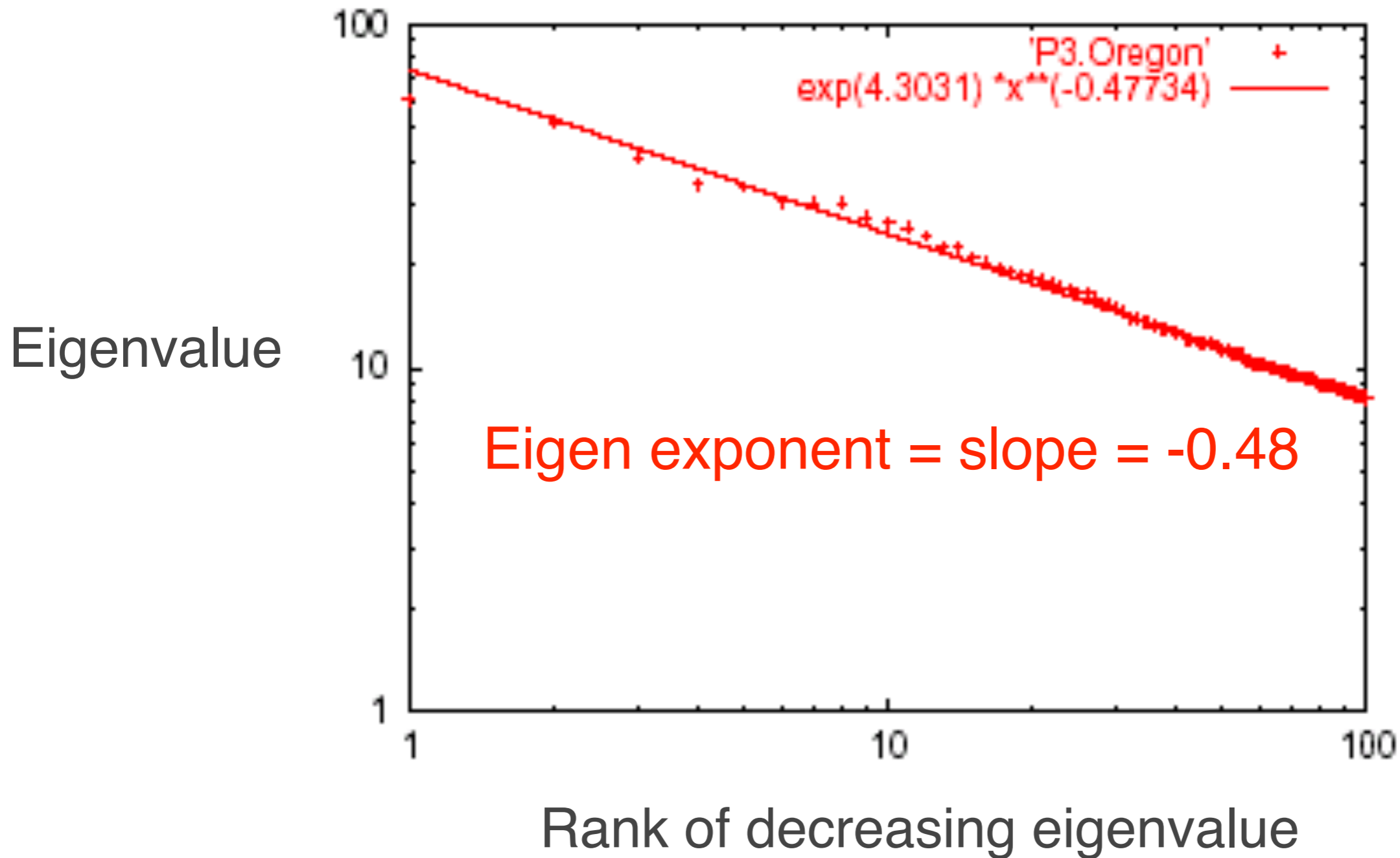
Q: Can we do that quickly?

A: Yes!

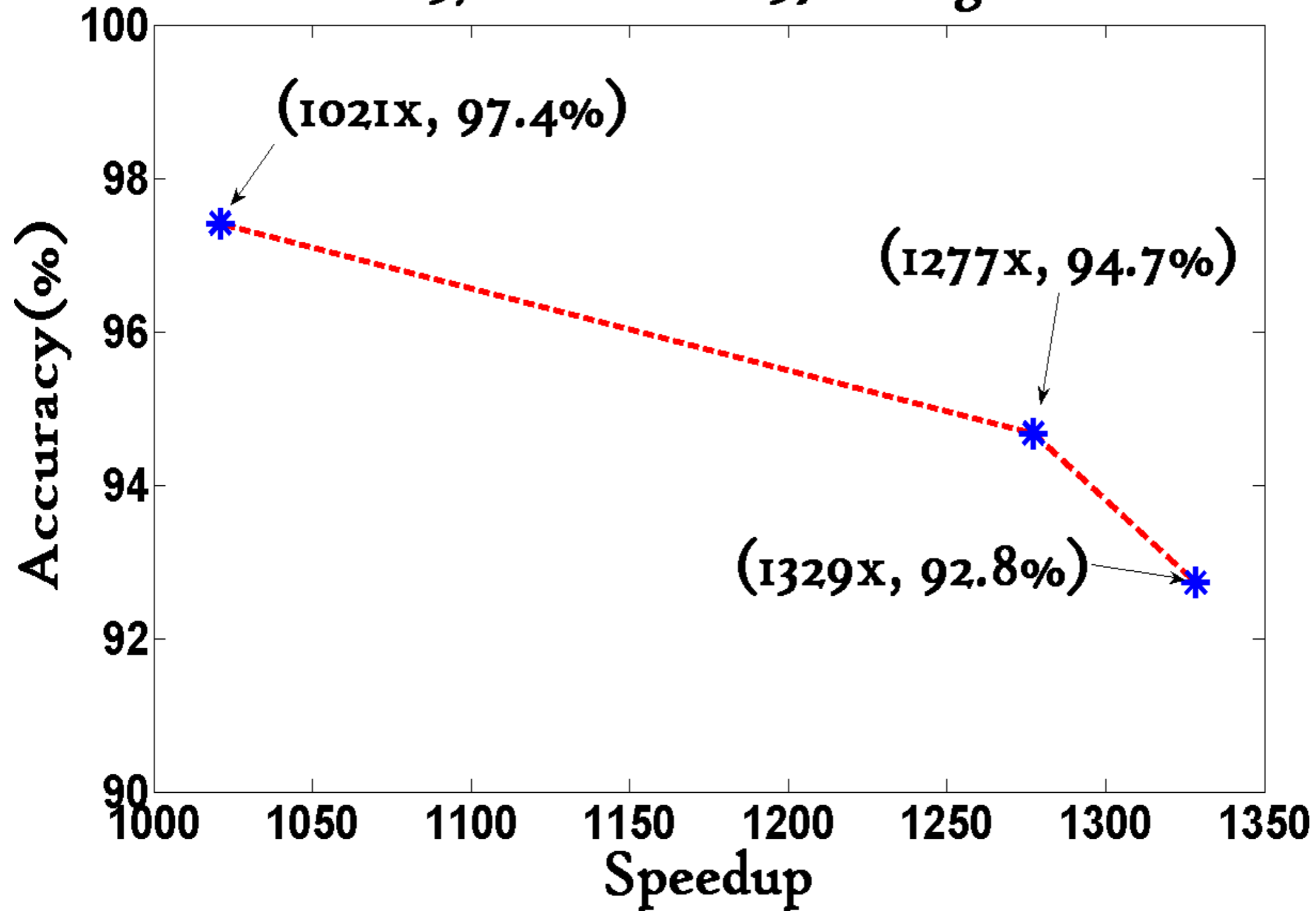
$$\#triangles = 1/6 \text{ Sum } ( \lambda_i^3 )$$

(and, because of skewness,  
we only need the top few eigenvalues!

# Power Law in Eigenvalues of Adjacency Matrix



Wikipedia graph 2006-Nov-04  
 $\approx 3,1\text{M}$  nodes  $\approx 37\text{M}$  edges



1000x+ speed-up, >90% accuracy

# More Centrality Measures...

- Degree
- Betweenness
- Closeness, by computing
  - Shortest paths
  - “**Proximity**” (usually via *random walks*) — **used successfully in a lot of applications**
- Eigenvector
- ...

# PageRank (Google)



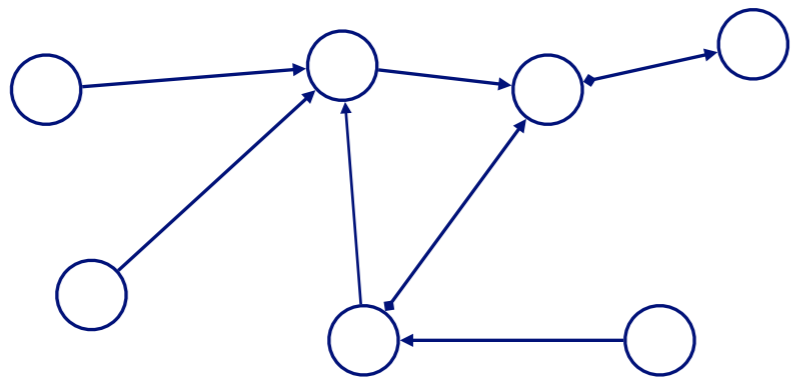
Larry Page

Sergey Brin

Brin, Sergey and Lawrence Page (1998).  
*Anatomy of a Large-Scale Hypertextual Web  
Search Engine*. 7th Intl World Wide Web Conf.

# PageRank: Problem

Given a directed graph, find its most interesting/central node



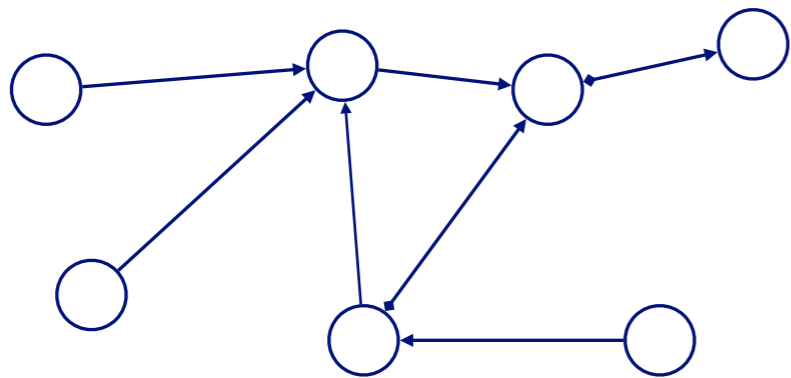
A node is important, if it is connected with important nodes (recursive, but OK!)

# PageRank: Solution

Given a directed graph, find its most interesting/central node

Proposed solution:

use **random walk**; spot most “popular” node  
(-> **steady state probability (ssp)**)



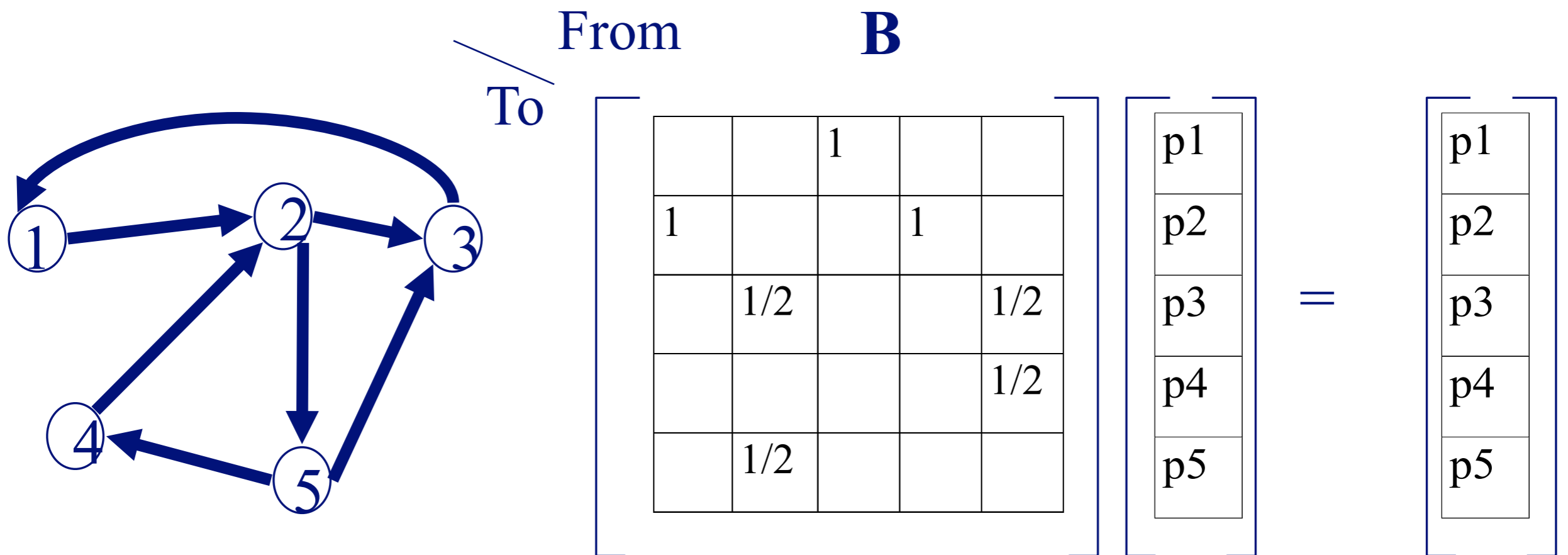
“state” = webpage

A node has high **ssp**,  
if it is connected  
with **high ssp** nodes  
(recursive, but OK!)



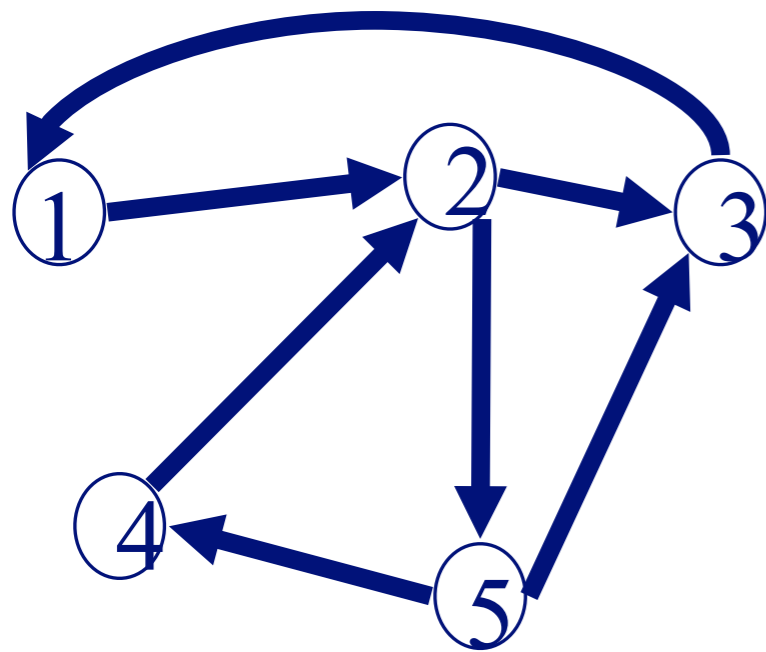
# (Simplified) PageRank

Let  $\mathbf{B}$  be the transition matrix:  
transposed, column-normalized



# (Simplified) PageRank

$$\mathbf{B} \mathbf{p} = \mathbf{p}$$



$$\mathbf{B} \mathbf{p} = \mathbf{p}$$

$$\begin{bmatrix} & & 1 & & \\ 1 & & & 1 & \\ & 1/2 & & & 1/2 \\ & & & & 1/2 \\ & 1/2 & & & \end{bmatrix} \begin{bmatrix} p1 \\ p2 \\ p3 \\ p4 \\ p5 \end{bmatrix} = \begin{bmatrix} p1 \\ p2 \\ p3 \\ p4 \\ p5 \end{bmatrix}$$

How to compute SSP:

<https://fenix.tecnico.ulisboa.pt/downloadFile/3779579688473/6.3.pdf>

<http://www.sosmath.com/matrix/markov/markov.html>

# (Simplified) PageRank

- $\mathbf{B} \mathbf{p} = \mathbf{1} * \mathbf{p}$
- Thus,  $\mathbf{p}$  is the **eigenvector** that corresponds to the highest eigenvalue ( $=1$ , since the matrix is column-normalized)
- Why does such a  $\mathbf{p}$  exist?
  - $\mathbf{p}$  exists if  $\mathbf{B}$  is  $n \times n$ , nonnegative, irreducible [Perron–Frobenius theorem]

# (Simplified) PageRank

- In short: imagine a particle randomly moving along the edges
- Compute its **steady-state probability (ssp)**

Full version of algorithm:

with **occasional random jumps**

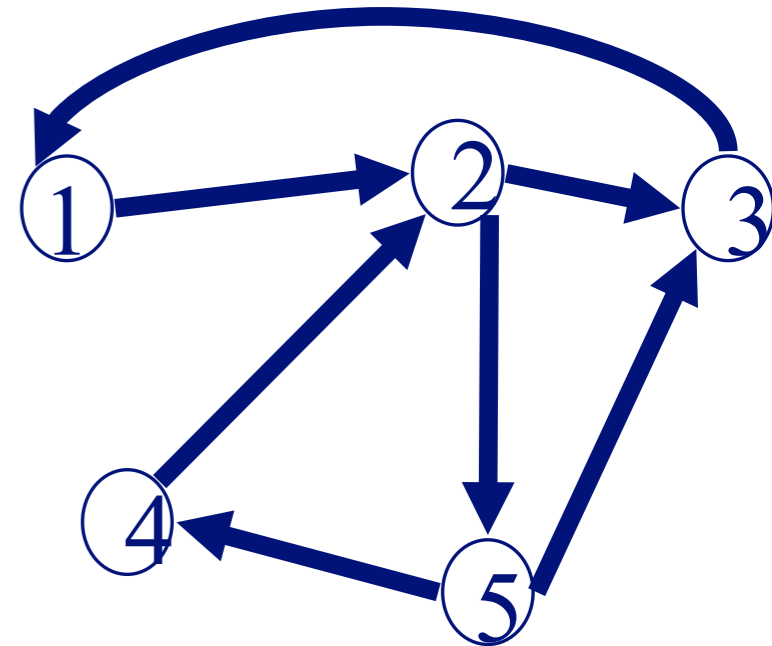
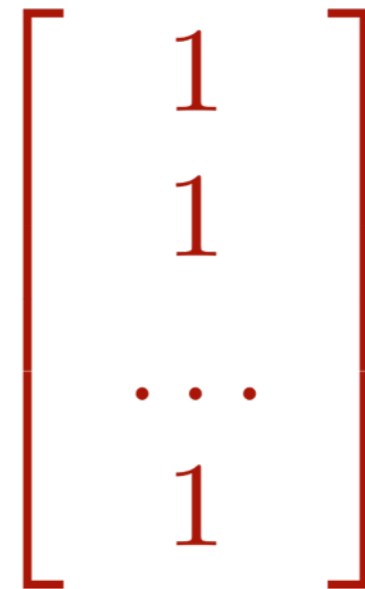
Why? To make the matrix irreducible

# Full Algorithm

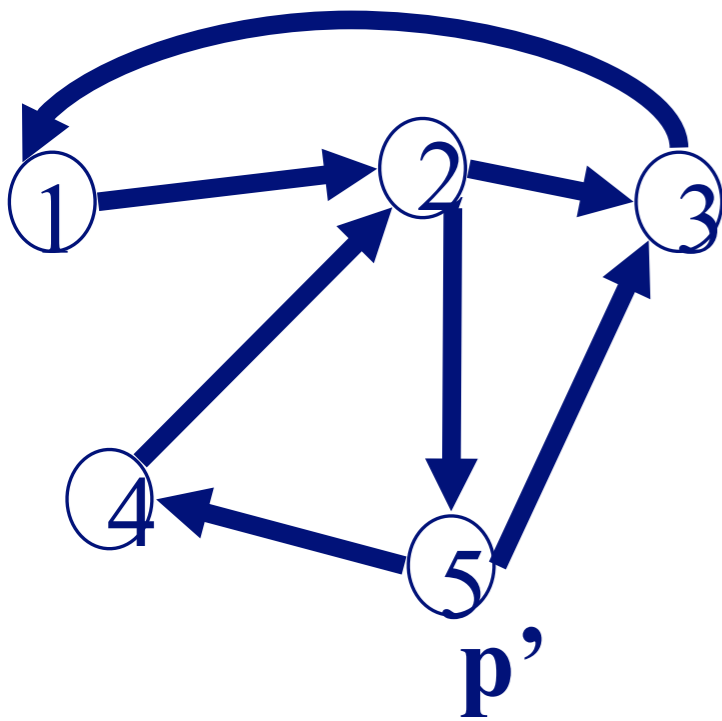
- With probability  $1-c$ , fly-out to a **random node**
- Then, we have

$$\mathbf{p} = c \mathbf{B} \mathbf{p} + (1-c)/n \mathbf{1} \Rightarrow$$

$$\mathbf{p} = (1-c)/n [\mathbf{I} - c \mathbf{B}]^{-1} \mathbf{1}$$



# How to compute PageRank for huge matrix?



Use the power iteration method

[http://en.wikipedia.org/wiki/Power\\_iteration](http://en.wikipedia.org/wiki/Power_iteration)

$$\mathbf{p} = c \mathbf{B} \mathbf{p} + (1-c)/n \mathbf{1}$$

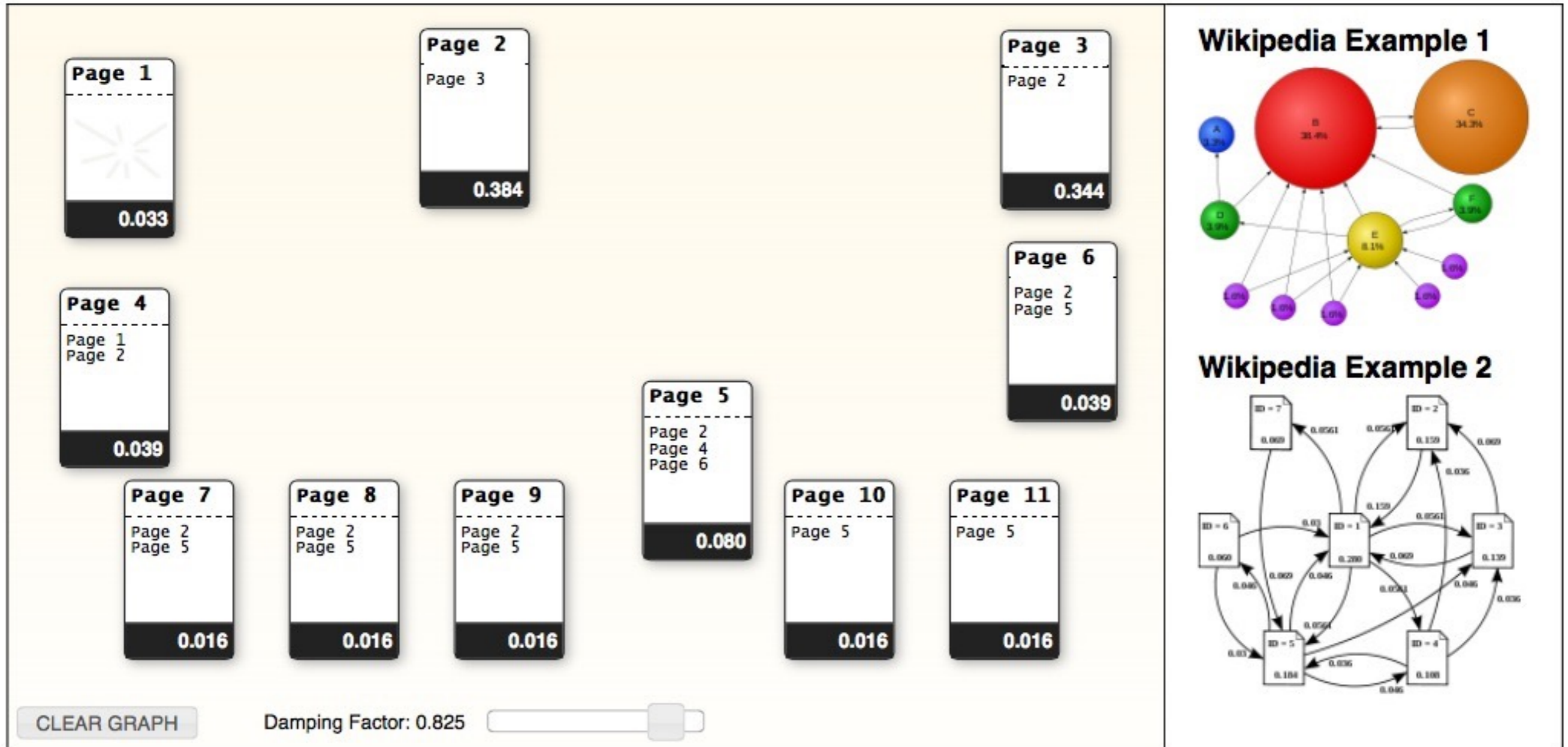
$\mathbf{B}$

$\mathbf{p}$

$$\begin{bmatrix} p1 \\ p2 \\ p3 \\ p4 \\ p5 \end{bmatrix} = c \begin{bmatrix} & & 1 & & \\ 1 & & & 1 & \\ & 1/2 & & & 1/2 \\ & & & & 1/2 \\ & 1/2 & & & \end{bmatrix} \begin{bmatrix} p1 \\ p2 \\ p3 \\ p4 \\ p5 \end{bmatrix} + (1-c)1/n \begin{bmatrix} 1 \\ 1 \\ \dots \\ 1 \end{bmatrix}$$

Can initialize this vector to any non-zero vector, e.g., all "1"s

# PageRank Explained with Javascript



<http://www.cs.duke.edu/csed/principles/pagerank/>

# PageRank for graphs (generally)

You can compute PageRank for **any graphs**

Should be in your algorithm “toolbox”

- Better than simple centrality measure (e.g., degree)
- Fast to compute for large graphs ( $O(E)$ )

But can be “misled” (Google Bomb)

- How?



# Personalized PageRank

Make one small variation of PageRank

- Intuition: not all pages are equal, some more relevant to a person's specific needs
- How?

# Personalized PageRank

With probability  $1-c$ , fly-out to ~~a random node~~ **some preferred nodes**

$$\mathbf{p}' = \mathbf{c} \mathbf{B} \mathbf{p} + (1-c) \mathbf{1}/n$$

The diagram illustrates the equation  $\mathbf{p}' = \mathbf{c} \mathbf{B} \mathbf{p} + (1-c) \mathbf{1}/n$ . On the left, a vertical vector  $\mathbf{p}'$  contains nodes p1, p2, p3, p4, and p5. In the middle, a 5x5 matrix  $\mathbf{B}$  is shown with the following values:
 

		1		
1			1	
	1/2			1/2
				1/2
	1/2			

 On the right, a vertical vector  $\mathbf{p}$  contains nodes p1, p2, p3, p4, and p5. A red vector with ones is shown next to it, with a blue arrow pointing to a green vector with 1s at the second and fifth positions. A small black arrow points to the vector  $\mathbf{p}$ .

Can initialize this vector to any non-zero vector, e.g., all "1"s

# Why learn Personalized PageRank?

## For recommendation

- If I like webpage A, what else do I like?
- If I bought product A, what other products would I also buy?

## Visualizing and interacting with large graphs

- Instead of visualizing every single nodes, visualize the **most important ones**

Very flexible — works on **any graph**

# Related “guilt-by-association” / diffusion techniques

- **Personalized PageRank**  
(= Random Walk with Restart)
- “Spreading activation” or “degree of interest”  
in Human-Computer Interaction (HCI)
- Belief Propagation  
(powerful inference algorithm, for fraud  
detection, image segmentation, error-  
correcting codes, etc.)

# Why are these algorithms popular?

- **Intuitive to interpret**  
uses “network effect”, homophily
- **Easy to implement**  
Math is relatively simple (mainly matrix-vector multiplication)
- **Fast**  
run time linear to #edges, or better
- **Probabilistic meaning**

# Human-In-The-Loop Graph Mining

**Apolo:**

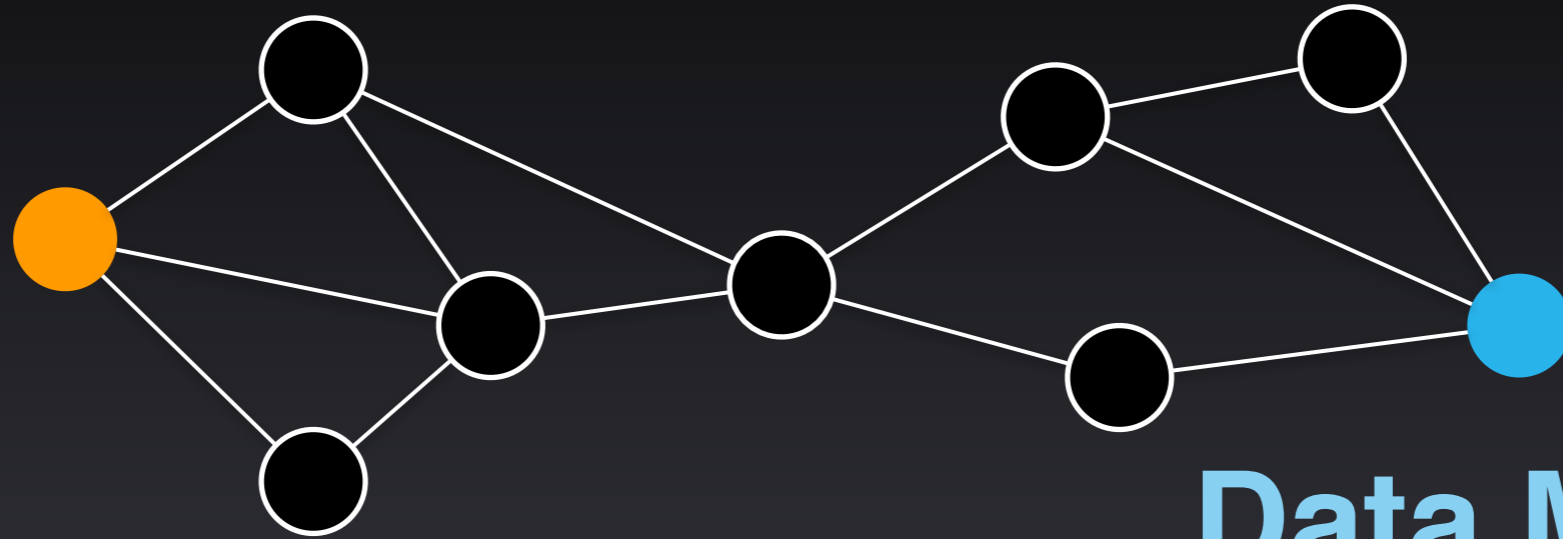
**Machine Learning + Visualization**

*CHI 2011*

Apolo: Making Sense of Large Network Data by Combining Rich User Interaction and Machine Learning

# Finding **More** Relevant Nodes

**HCI**  
Paper

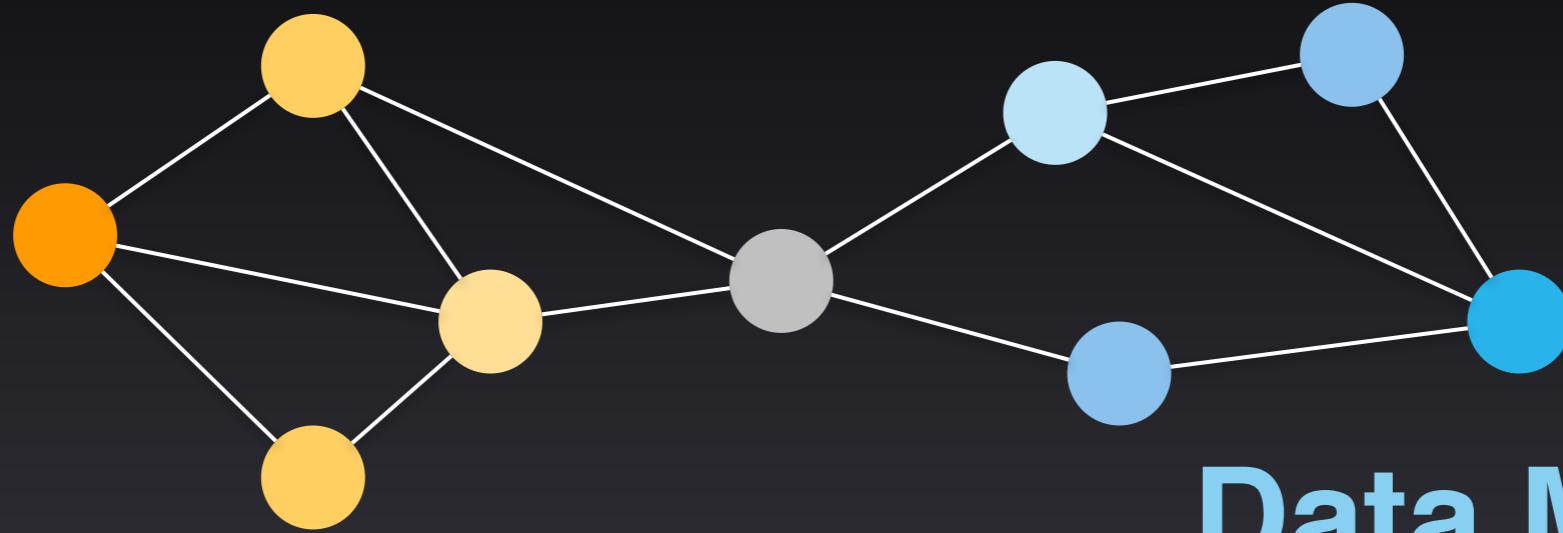


**Data Mining**  
Paper

Citation network

# Finding **More** Relevant Nodes

**HCI**  
Paper

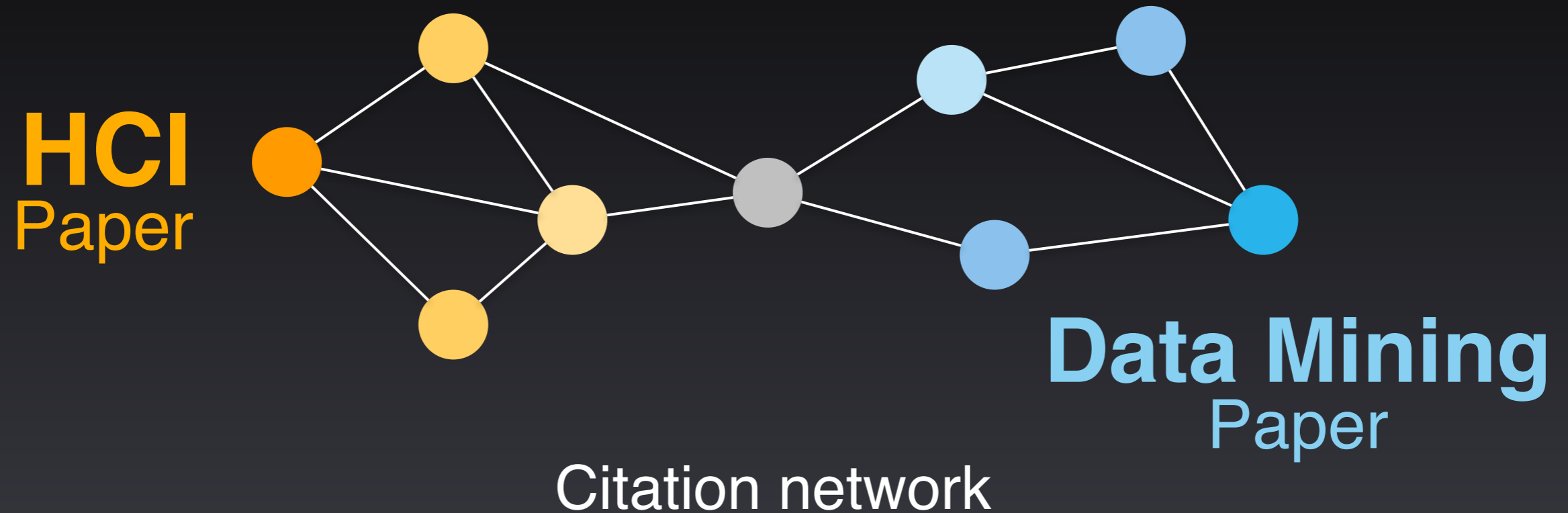


**Data Mining**  
Paper

Citation network



# Finding **More** Relevant Nodes

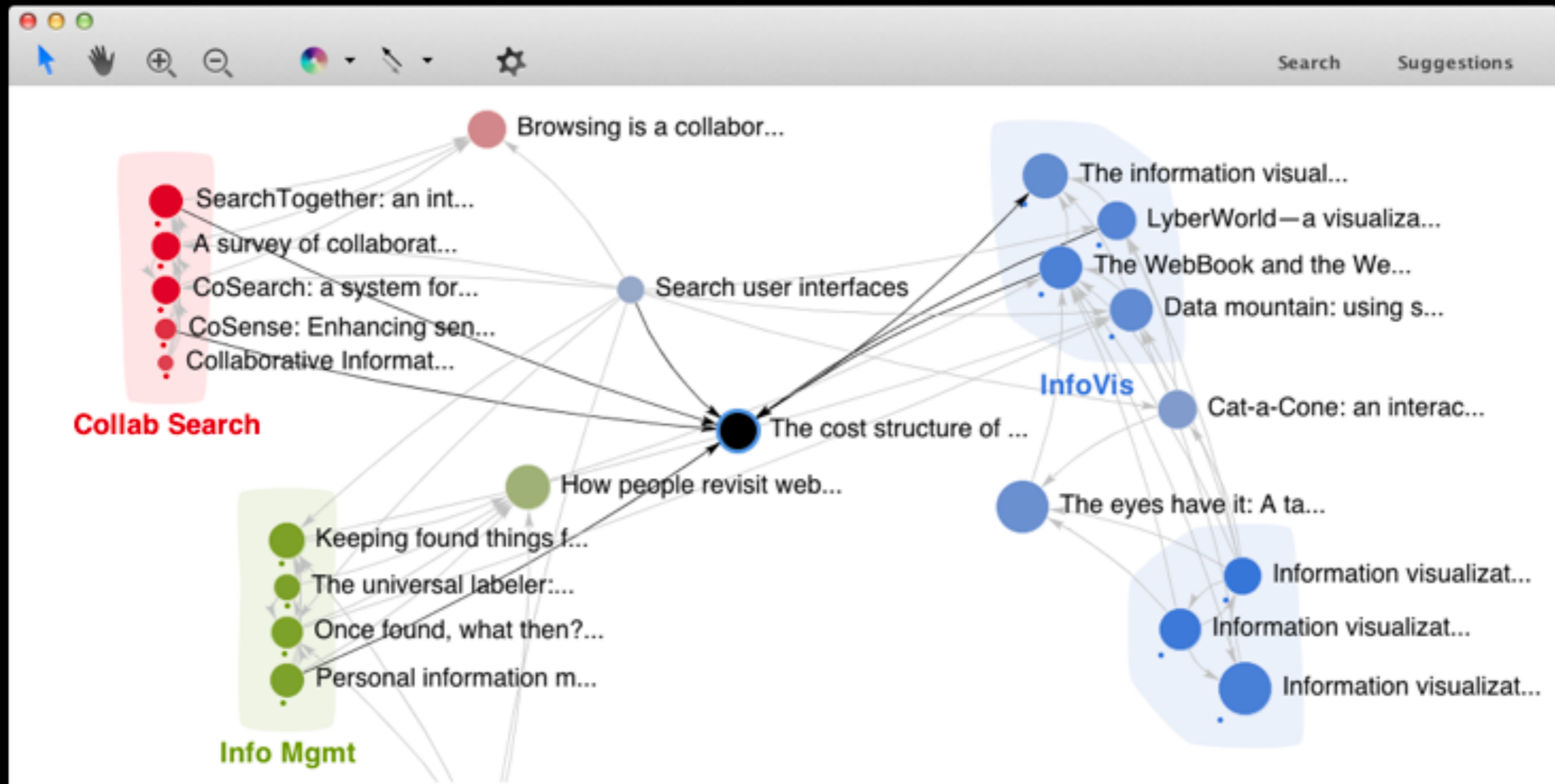


Apolo uses **guilt-by-association**  
(Belief Propagation, similar to personalized PageRank)

# Demo: Mapping the Sensemaking Literature

Nodes: 80k papers from Google Scholar (node size: #citation)

Edges: 150k citations



**The cost structure of sensemaking**

Russell, D.M. and Stefik, M.J. and Pirolli, P. and Card, S.K.

245 citations 8 versions

PDF 1993

For The cost structure of sensemaking

<b>The information visualizer, an inf...</b>	<b>1991</b>
Card, S.K. and Robertson, G.G. and Macki...	532
<b>The WebBook and the Web Forag...</b>	<b>1996</b>
Card, S.K. and Robertson, G.G. and York, W.	403
<b>LyberWorld—a visualization user...</b>	<b>1994</b>
Hemmje, M. and Kunkel, C. and Willett, A.	223
<b>The structure of the information...</b>	<b>1997</b>
Card, S.K. and Mackinlay, J.	198
<b>Information visualization</b>	<b>2009</b>
Card, S. and Mackinlay, JD and Shneiderm...	180
<b>"I'll get that off the audio": a cas...</b>	<b>1997</b>
Moran, T.P. and Palen, L. and Harrison, S...	143
<b>An organic user interface for sear...</b>	<b>1995</b>
Mackinlay, J.D. and Rao, R. and Card, S.K.	123
<b>Using a landscape metaphor to re...</b>	<b>1993</b>
Chalmers, M.	122
<b>Personal information management</b>	<b>2007</b>
Jones, W.P. and Teevan, J.	109
<b>SearchTogether: an interface for c...</b>	<b>2007</b>
Morris, M.R. and Horvitz, E.	108
<b>Information foraging theory: Ada...</b>	<b>2007</b>
Pirolli, P.	107
<b>Investigating behavioral variabilit...</b>	<b>2007</b>
White, R.W. and Drucker, S.M.	79
<b>Jigsaw: Supporting investigative...</b>	<b>2008</b>
Stasko, J. and Görg, C. and Liu, Z.	71
<b>The cost-of-knowledge character...</b>	<b>1994</b>
Card, S.K. and Pirolli, P. and Mackinlay, J.D.	54
<b>Collaborative conceptual design:...</b>	<b>1996</b>
Potts, C. and Catledge, L.	45

The cost structure of sen...

**The cost structure of sensemaking** PDF 1993  
*Russell, D.M. and Stefik, M.J. and Pirolli, P. and Card, S.K.*  
 245 citations 8 versions

For The cost structure of sensemaking

<b>The information visualizer, an inf...</b>	<b>1991</b>
Card, S.K. and Robertson, G.G. and Macki...	532
<b>The WebBook and the Web Forag...</b>	<b>1996</b>
Card, S.K. and Robertson, G.G. and York, W.	403
<b>LyberWorld—a visualization user...</b>	<b>1994</b>
Hemmje, M. and Kunkel, C. and Willett, A.	223
<b>The structure of the information...</b>	<b>1997</b>
Card, S.K. and Mackinlay, J.	198
<b>Information visualization</b>	<b>2009</b>
Card, S. and Mackinlay, JD and Shneiderm...	180
<b>"I'll get that off the audio": a cas...</b>	<b>1997</b>
Moran, T.P. and Palen, L. and Harrison, S...	143
<b>An organic user interface for sear...</b>	<b>1995</b>
Mackinlay, J.D. and Rao, R. and Card, S.K.	123
<b>Using a landscape metaphor to re...</b>	<b>1993</b>
Chalmers, M.	122
<b>Personal information management</b>	<b>2007</b>
Jones, W.P. and Teevan, J.	109
<b>SearchTogether: an interface for c...</b>	<b>2007</b>
Morris, M.R. and Horvitz, E.	108
<b>Information foraging theory: Ada...</b>	<b>2007</b>
Pirolli, P.	107
<b>Investigating behavioral variabilit...</b>	<b>2007</b>
White, R.W. and Drucker, S.M.	79
<b>Jigsaw: Supporting investigative...</b>	<b>2008</b>
Stasko, J. and Görg, C. and Liu, Z.	71
<b>The cost-of-knowledge character...</b>	<b>1994</b>
Card, S.K. and Pirolli, P. and Mackinlay, J.D.	54
<b>Collaborative conceptual design:...</b>	<b>1996</b>
Potts, C. and Catledge, L.	45

The cost structure of sen...

**The cost structure of sensemaking** PDF 1993

*Russell, D.M. and Stefik, M.J. and Pirolli, P. and Card, S.K.*

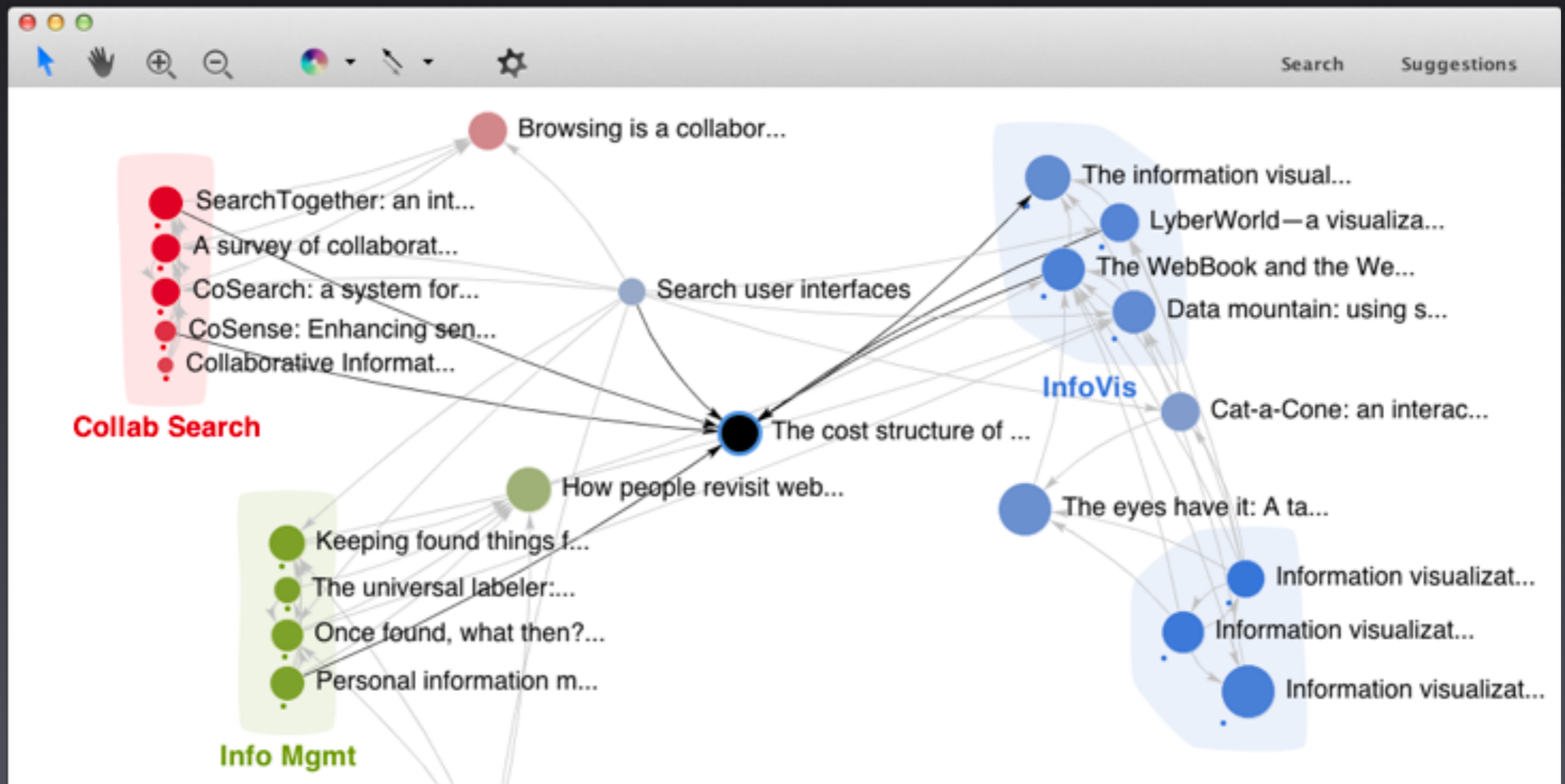
245 citations 8 versions

# Key Ideas (Recap)



Specify **exemplars**

Find **other** relevant nodes (BP)



# Apolo's Contributions

## 1 Human + Machine

It was like having a  
partnership with the machine.



Apolo User

## 2 Personalized Landscape

# Apolo 2009

The screenshot displays the Apolo 2009 interface, which is a search engine for research papers. At the top, there are navigation buttons for "Cluster Data" and "Add Group", and a "Recommendations" slider. The main area shows a list of search results, with several clusters highlighted in different colors:

- End User Programming (Green):** This cluster includes papers such as "End users creating effective softw...", "End user software engineering: chi...", and "Invited research overview: end-us...". The author "Brad A. Myers" is prominently displayed.
- Not Interested (Blue):** This cluster includes papers like "Automatically generating user inte...", "Decision-Theoretic User Interface...", and "Daniel S. Weld".
- Text Entry (Light Blue):** This cluster includes papers such as "In-stroke word completion.", "Integrating isometric joysticks into...", and "Eyes on the road, hands on the whe...".
- Interface Generation (Orange):** This cluster includes papers like "Huddle: automatically generating i...", "UNIFORM: automatically generatin...", and "Demonstrating the viability of auto...".
- Brad (Yellow):** This cluster is centered around the author "Brad A. Myers" and includes papers like "The garnet user interface developm...", "Using HCI Techniques to Design a M...", and "Creating charts by demonstration."

At the bottom of each cluster, there is a "Show:" dropdown menu, with "All" selected for the green cluster and "Papers" selected for the light blue cluster.

# Apolo 2010

Shiftr

Date Save/Load Export

Search  103 matches

all title authors

Title	Cites	Authors	Year
The cost structure of sensemaking	188	Russell, D	1993
Table lens as a tool for making sense	37	Pirolli, P.	1996
Sensemaking of evolving web sites	22	Chi, E.H.	1999
Sources of structure in sensemaking	19	Qu, Y. an	2005
A sensemaking-supporting informa	11	Qu, Y.	2003

PIM Collab search InfoVis **Sensemaking**

Title	Cit...	Authors	Year
SenseMaker: an information-explor	155	Baldonad	1997
Sources of structure in sensemakin	19	Qu, Y. an	2005
The cost structure of sensemaking	188	Russell, D	1993
Inferring web communities from li	663	Gibson, D	1998
Sensemaking for Topic Compreher	0	Ryder, B.	0
A sensemaking-supporting inform	11	Qu, Y.	2003
Model-driven formative evaluation	6	Qu, Y. an	2008
The digital library integrated task e	85	Cousins,	1997
CiteSense: supporting sensemaking	0	Zhang, X.	2008
An informal information-seeking e	53	Hendry, I	1997
An empirical evaluation of user inti	35	Amento,	1999
The effectiveness of automatically	19	Gonf{c}{c	2004
The microstructures of social taggi	3	Fu, W.T.	2008
Sensemaking: Bringing theories an	0	Sharma, I	2006
Data manipulation services in the I	7	Asdooria	1998
Considerations for information env	78	Fumas, G	1998

17 nodes selected

**1993**  
**The cost structure of sensemaking**  
 Russell, D.M., Stefik, M.J., Pirolli, P., Card, S.K.  
 Cited by 188

booktitle Proceedings of the INTERACT'93 and CHI'93 con

Change in representations

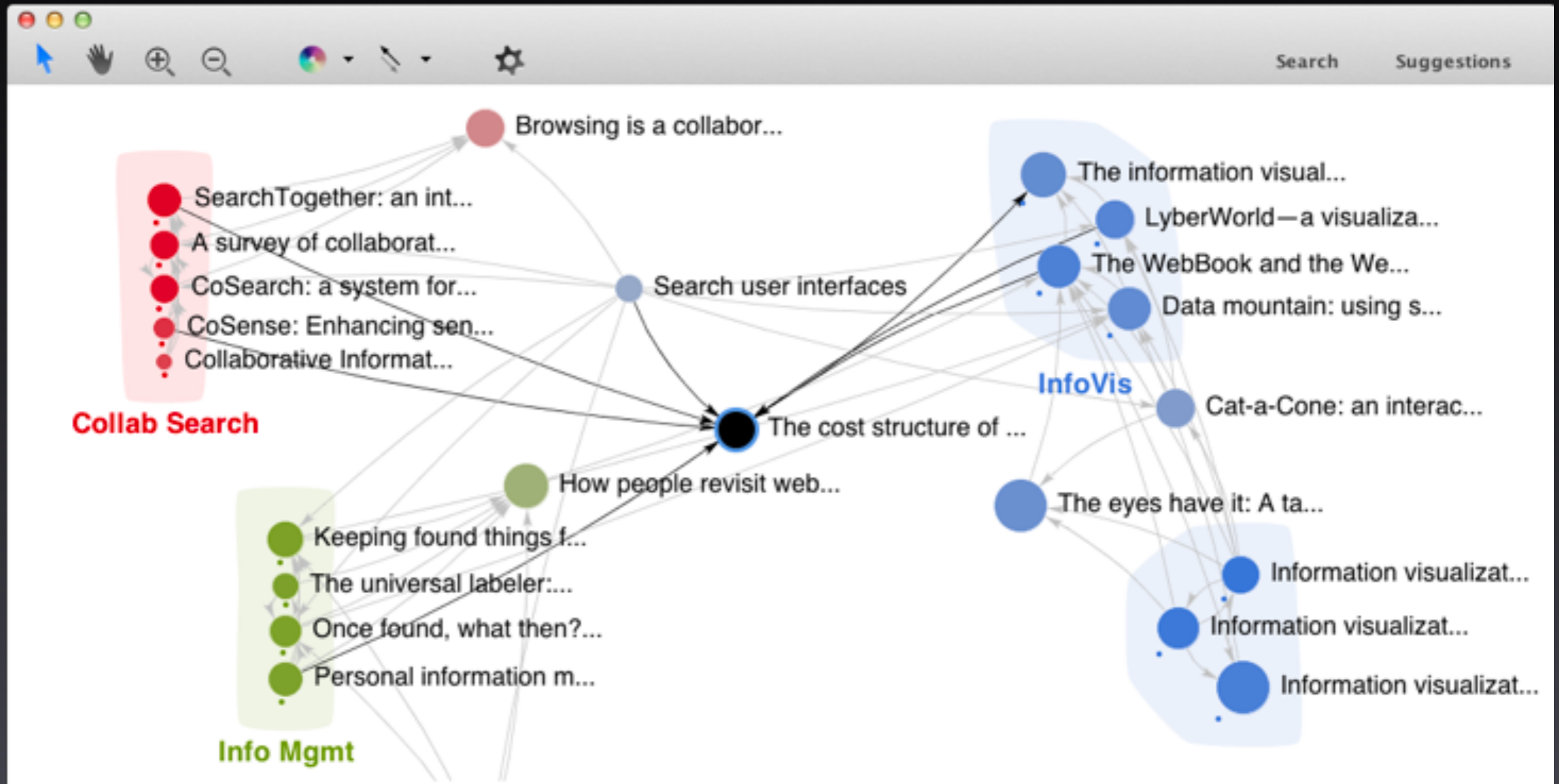
**Change in Representations**

Iteration	Red	Blue	Green	Orange
0	1	1	1	1
1	1	2	1	3
2	1	2	4	6
3	1	3	4	7
4	1	3	8	8
5	2	5	8	8
6	2	5	8	9
7	2	5	8	9



# Apolo 2011

22,000 lines of code. Java 1.6. Swing.  
Uses SQLite3 to store graph on disk



**The cost structure of sensemaking**

Russell, D.M. and Stefik, M.J. and Pirolli, P. and Card, S.K.

245 citations 8 versions

PDF 1993

# User Study

Used **citation network**

**Task:** Find related papers for **2 sections** in a **survey paper on *user interface***

- **Model-based** generation of UI
- Rapid **prototyping** tools



**Past, Present and Future of  
User Interface Software Tools**

**Brad Myers, Scott E. Hudson, and Randy Pausch**

Human Computer Interaction Institute  
School of Computer Science  
Carnegie Mellon University  
Pittsburgh, PA 15213-3891

**Apolo**



**Google Scholar**



**Between subjects design**

**Participants: grad student or research staff**

# Apolo

# Google Scholar



“Model-based”



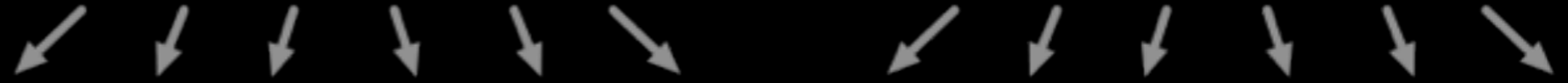
“Prototyping”



*10 papers for each section*

# Apolo

# Google Scholar



“Model-based”

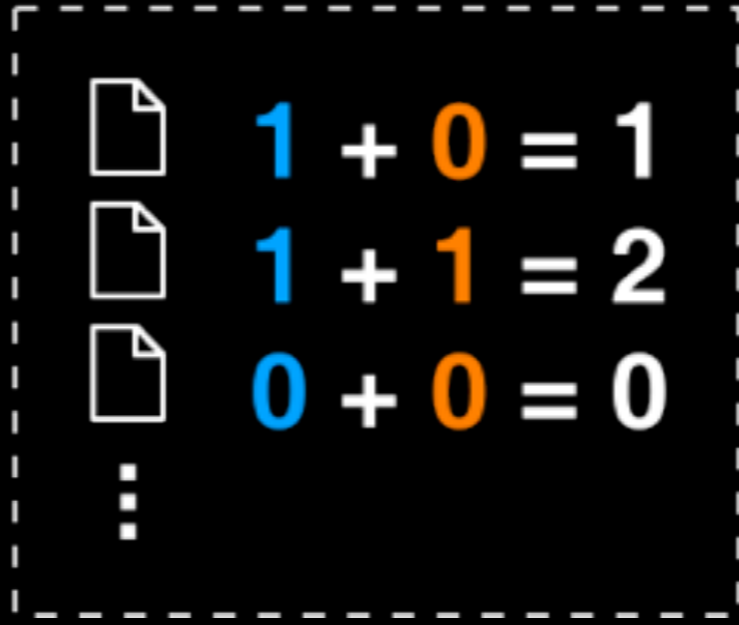


“Prototyping”

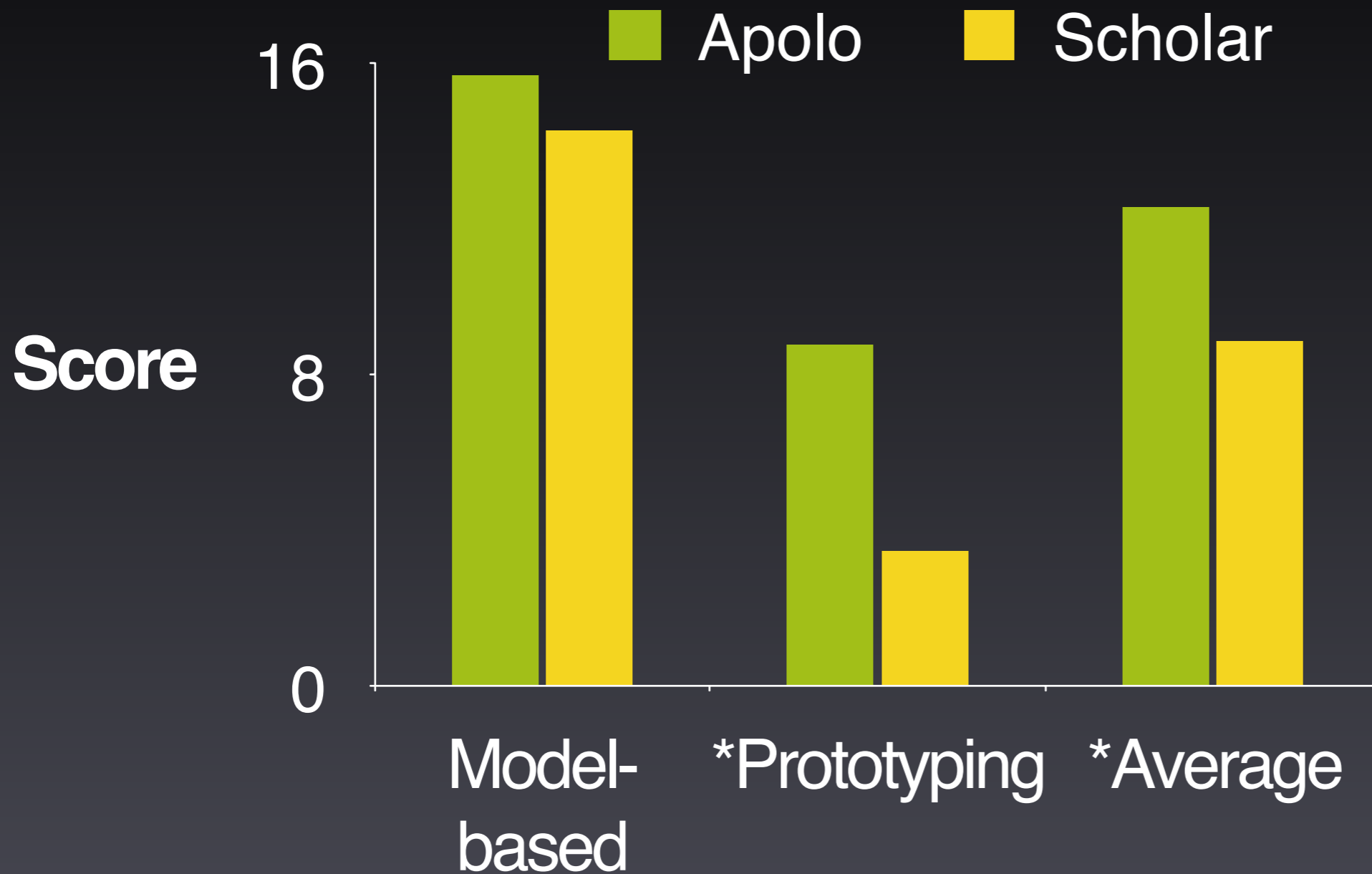
*10 papers for each section*



Expert judges rated papers



# Judges' Scores



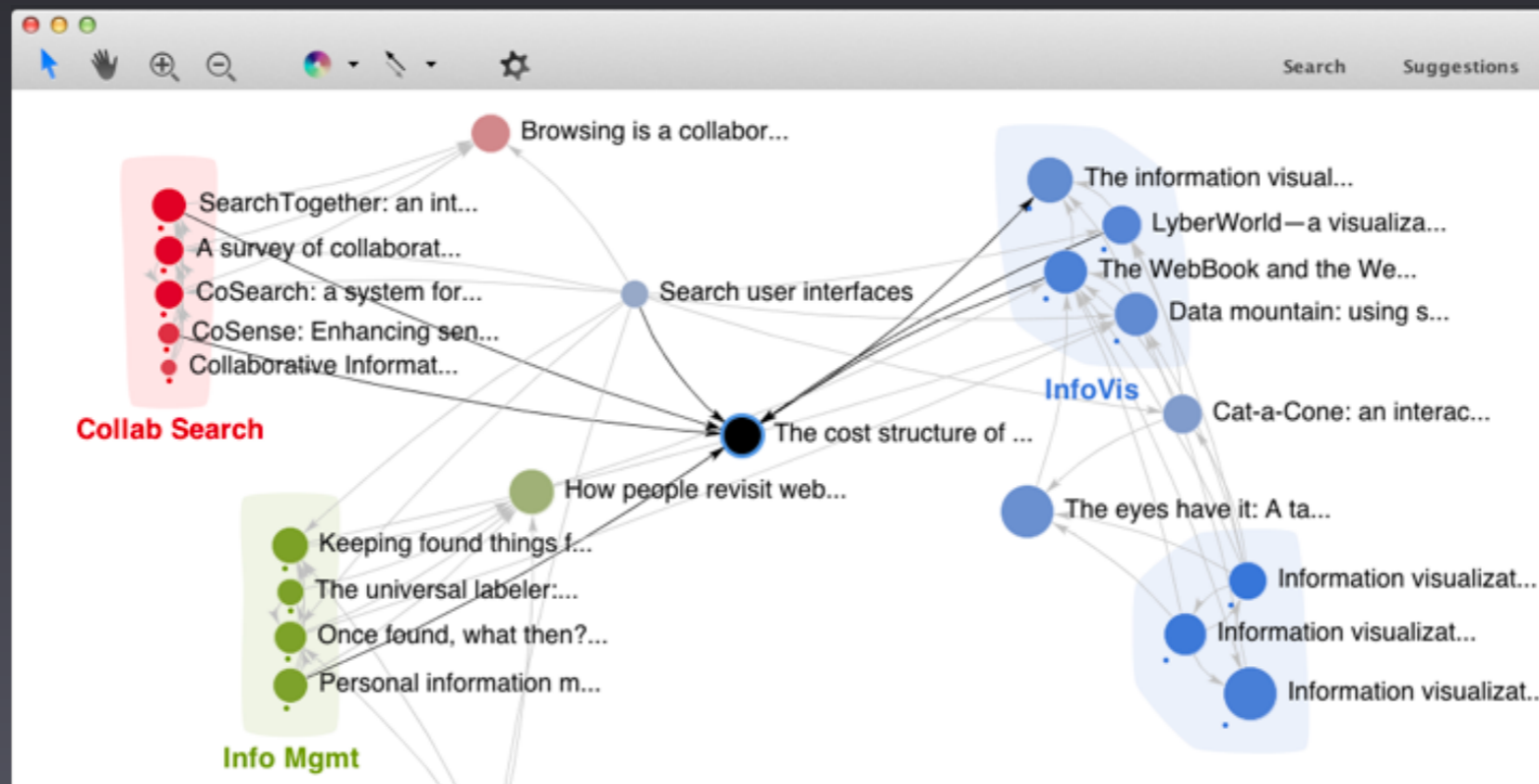
↑  
Higher is better.  
**Apolo** wins.

\* Statistically significant, by *two-tailed t test*,  $p < 0.05$

# Apolo: Recap

A **mixed-initiative** approach for **exploring** and creating personalized **landscape** for large network data

Apolo = ML + Visualization + Interaction



# Practitioners' guide to building (interactive) applications

Think about scalability early

- Identify candidate scalable algorithms early on

Use **iterative** design approach, as in Apollo and industry

- Why? It's hard to get it right the first time
- **Create prototype, evaluate, modify prototype, evaluate, ...**
- Quick evaluation helps you identify **important fixes early — save you a lot of time overall**



# Practitioners' guide to building (interactive) applications

What kinds of **prototypes**?

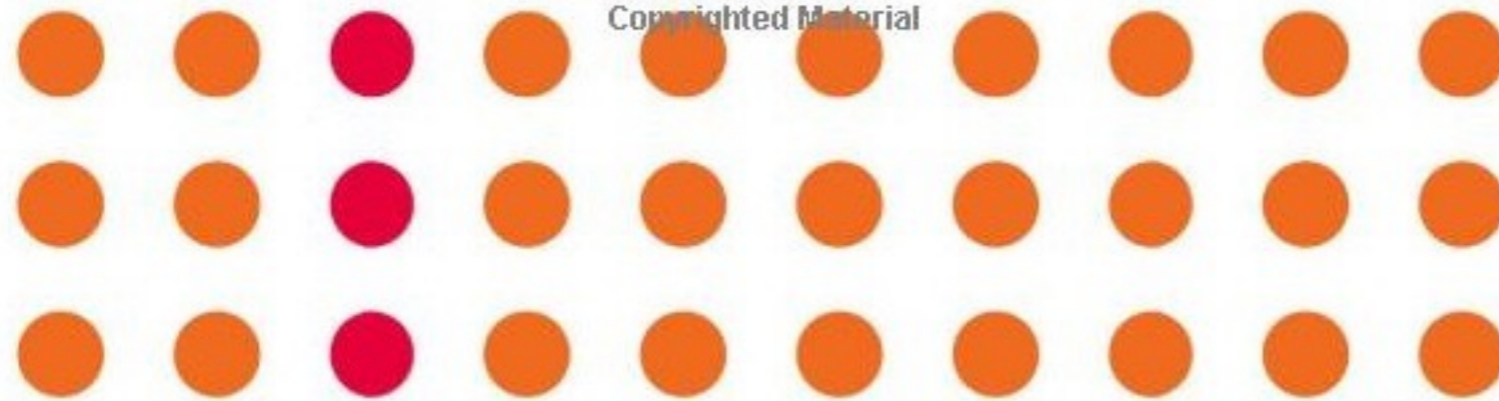
- Paper prototype, lo-fi prototype, high-fi prototype

Important to involve **REAL users** as early as possible

- Recruit your friends to try your tools
- Lab study (controlled, as in Apollo)
- Longitudinal study (usage over months)
- Deploy it and see the world's reaction!
- To learn more:
  - CS 6750 Human-Computer Interaction
  - CS 6455 User Interface Design and Evaluation

# If you want to know more about people...

<http://amzn.com/0321767535>



# 100 THINGS

EVERY DESIGNER NEEDS TO KNOW ABOUT **PEOPLE**

SUSAN M. WEINSCHENK, Ph.D.

