

PageRank (Google)



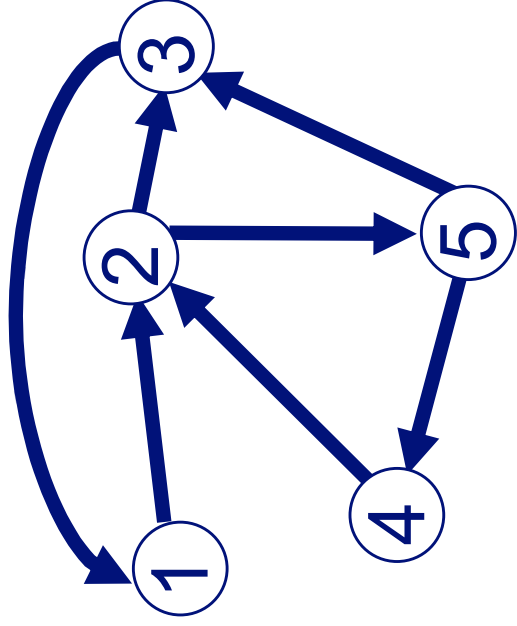
Larry Page

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

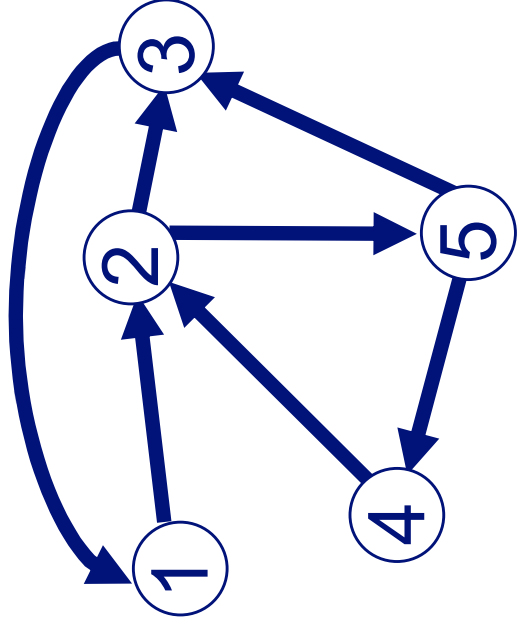
Given a directed graph, find its most interesting/central node



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

PageRank: Solution

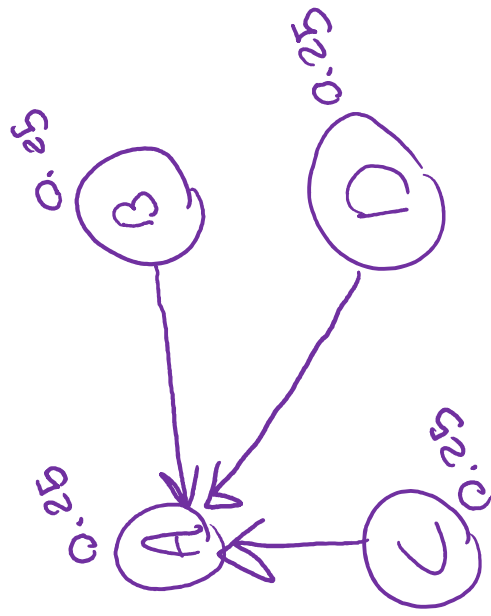
Given a directed graph, find its most interesting/central node
Proposed solution: use **random walk**; most “popular” nodes are the ones with highest **steady state probability (ssp)**



“**state**” = webpage

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

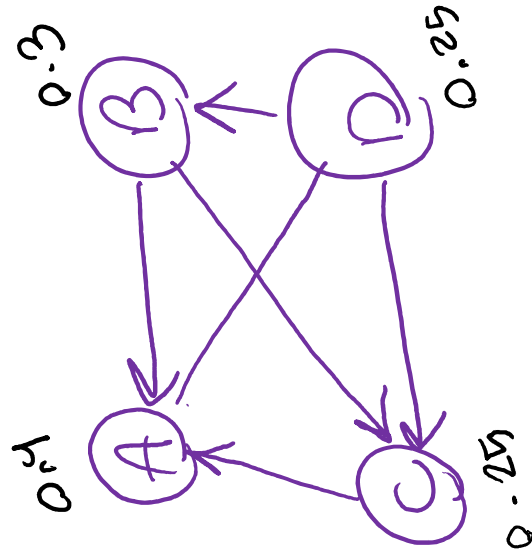
$$\frac{1}{4} = 0.25$$



$$PR(A) = \cancel{PR(B)} + PR(C) + PR(D)$$

$$= 0.25 + 0.25 + 0.25 = 0.75$$

$$\frac{0.75}{3} = 0.25$$

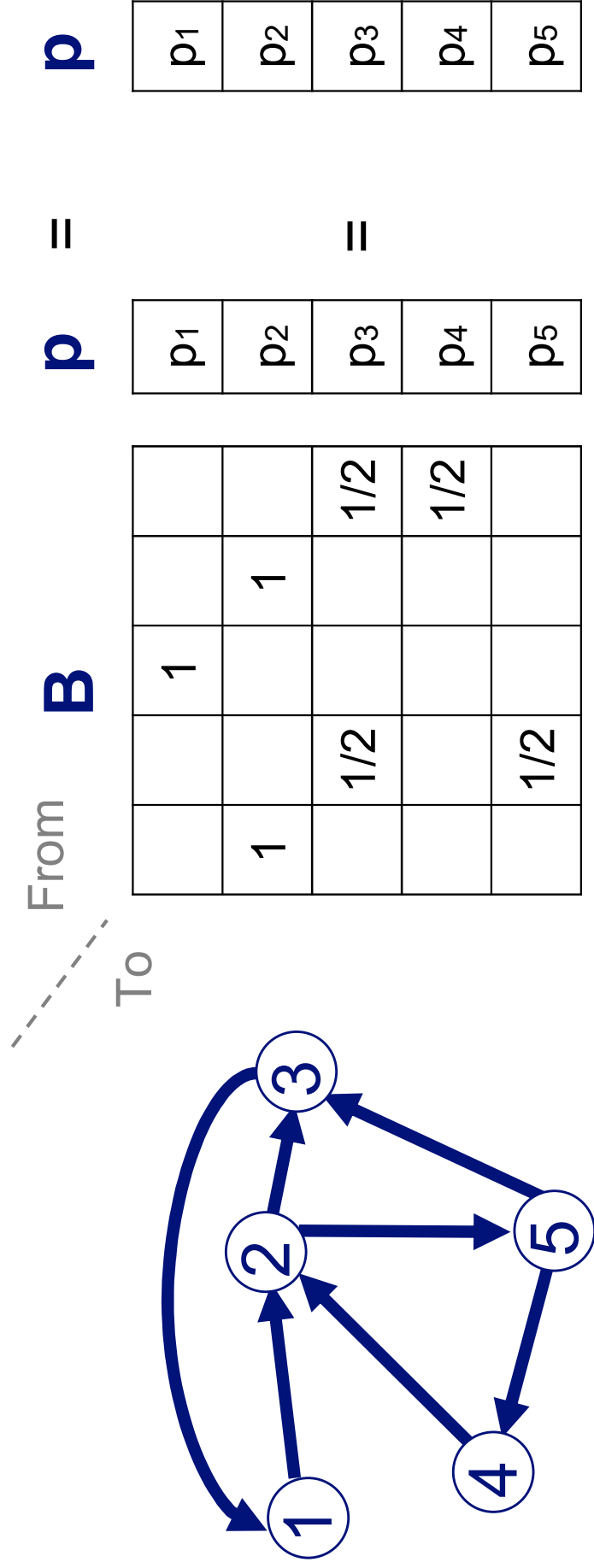


$$PR(A) = \frac{0.25}{2} + 0.25 + \frac{PR(B)}{2} + \frac{PR(C)}{3} + \frac{PR(D)}{3}$$

$$\cancel{PR(A)} = \frac{0.25}{2} + \frac{PR(B)}{2} + \frac{PR(C)}{3} + \frac{PR(D)}{3}$$

(Simplified) PageRank

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



How to compute SSP:

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

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

(Simplified) PageRank

$$A \tilde{x} = \tilde{x}$$

$$B \mathbf{p} = \mathbf{1}$$

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 person **randomly moving** along the edges/links
- A node's PageRank score is the **steady-state probability (ssp)** of finding the person at that node

Full version of algorithm:

With **occasional random jumps to any nodes**

Why? To make the matrix **irreducible**.

Irreducible = from any state (node), there's **non-zero probability to reach any other state** (node)

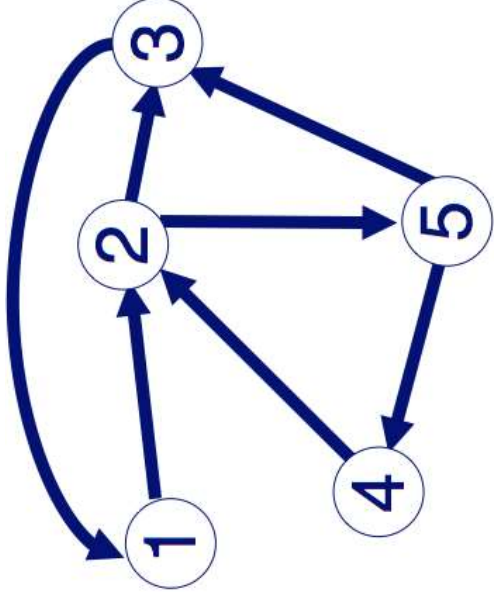
Full Algorithm

$$p = \beta p$$

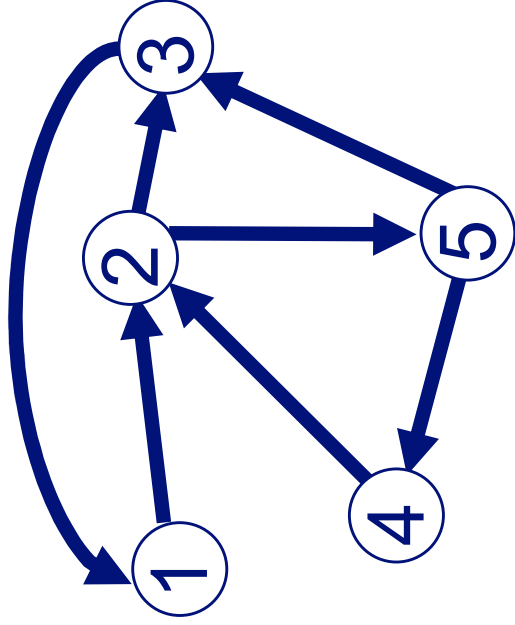
With probability $1-c$, fly-out to a **random node**

Then, we have

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



How to compute PageRank for huge matrix?



Use the power iteration method

http://en.wikipedia.org/wiki/Power_iteration

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

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

$$\mathbf{p}' = \begin{bmatrix} p1 \\ p2 \\ p3 \\ p4 \\ p5 \end{bmatrix}$$

= C

$$\mathbf{B} = \begin{bmatrix} & & & & \\ 1 & & & & \\ & 1 & & & \\ & & 1/2 & & \\ & & & 1/2 & \\ & & & & 1/2 \end{bmatrix}$$

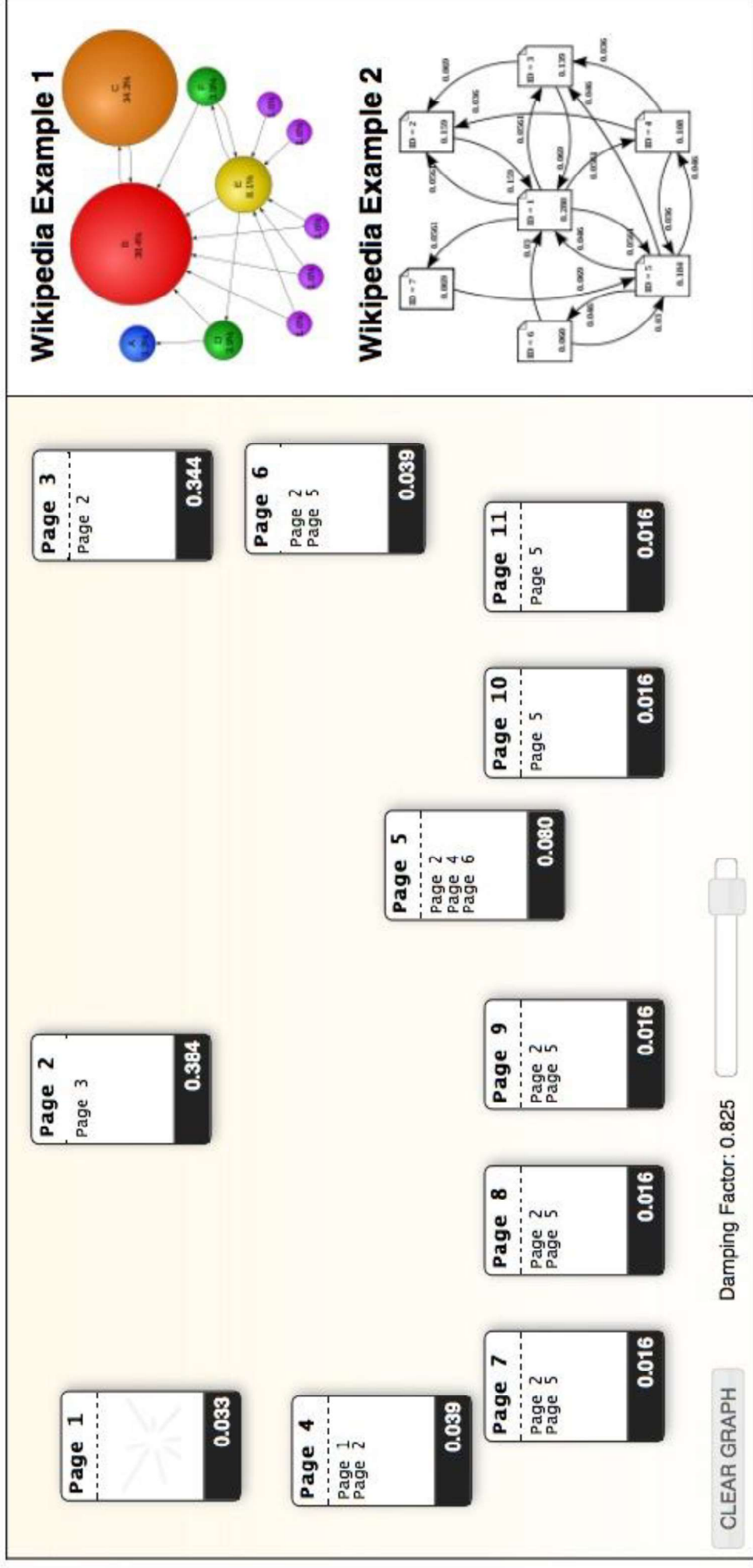
$$\mathbf{p} = \begin{bmatrix} p1 \\ p2 \\ p3 \\ p4 \\ p5 \end{bmatrix}$$

$$+ \frac{(1-c)}{n}$$

$$\begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix}$$

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

PageRank Explained with Javascript



Also great for checking the correctness of your PageRank Implementation.

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

PageRank for graphs (generally)

You can run PageRank on **any graphs**

- All you need are the **graph edges!**

Should be in your algorithm “toolbox”

- Better than degree centrality
- Fast to compute for large graphs, runtime linear in the number of edges, $O(E)$

But can be “misled” (Google Bomb)

- How?

Personalized PageRank

Intuition: not all pages are equal, some more relevant to some people

Goal: rank pages in a way that those more relevant to you will be ranked higher

How? Make just **one** small change to PageRank

Personalized PageRank

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

$$\mathbf{p}' = c \mathbf{B} \mathbf{p} + \frac{1-c}{N} \mathbf{1}$$

p'_1
p'_2
p'_3
p'_4
p'_5

$= 0.8$

Default value for c

		1			p_1
1			1		p_2
				$1/2$	p_3
				$1/2$	p_4
				$1/2$	p_5

$+ \frac{0.2}{5}$

1
1
1
1
1

0
1
0
0
1

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?

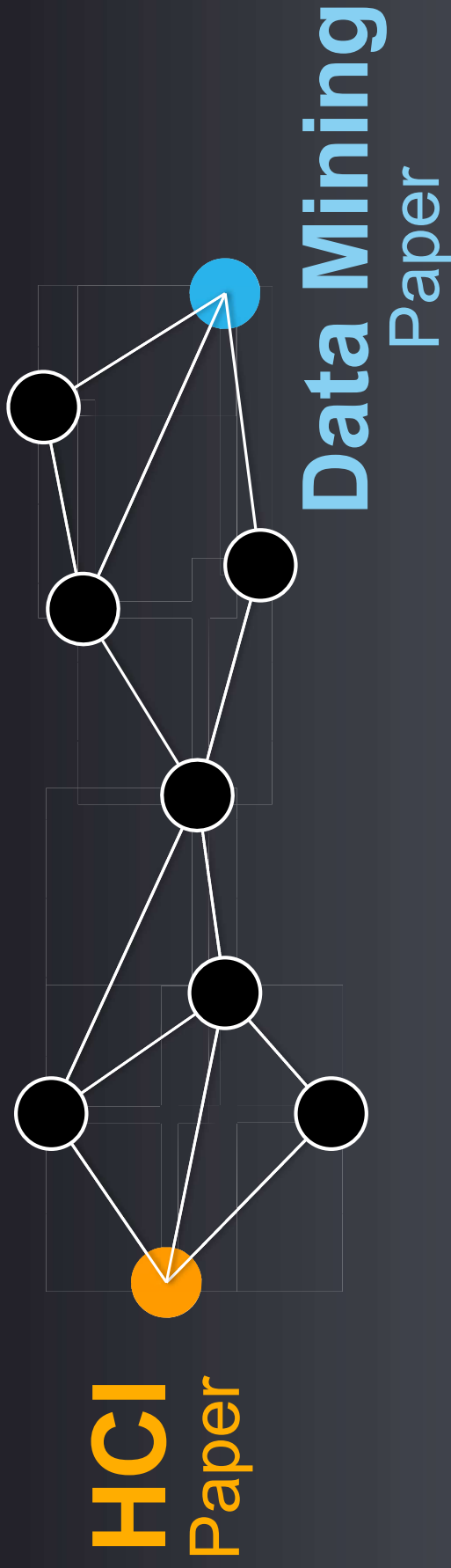
- **Intuitive to interpret**
uses “network effect”, homophily
- **Easy to implement**
math is relatively simple (mainly matrix-vector multiplication)
- **Fast**
run time linear to #edges, or better
- **Probabilistic meaning**

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

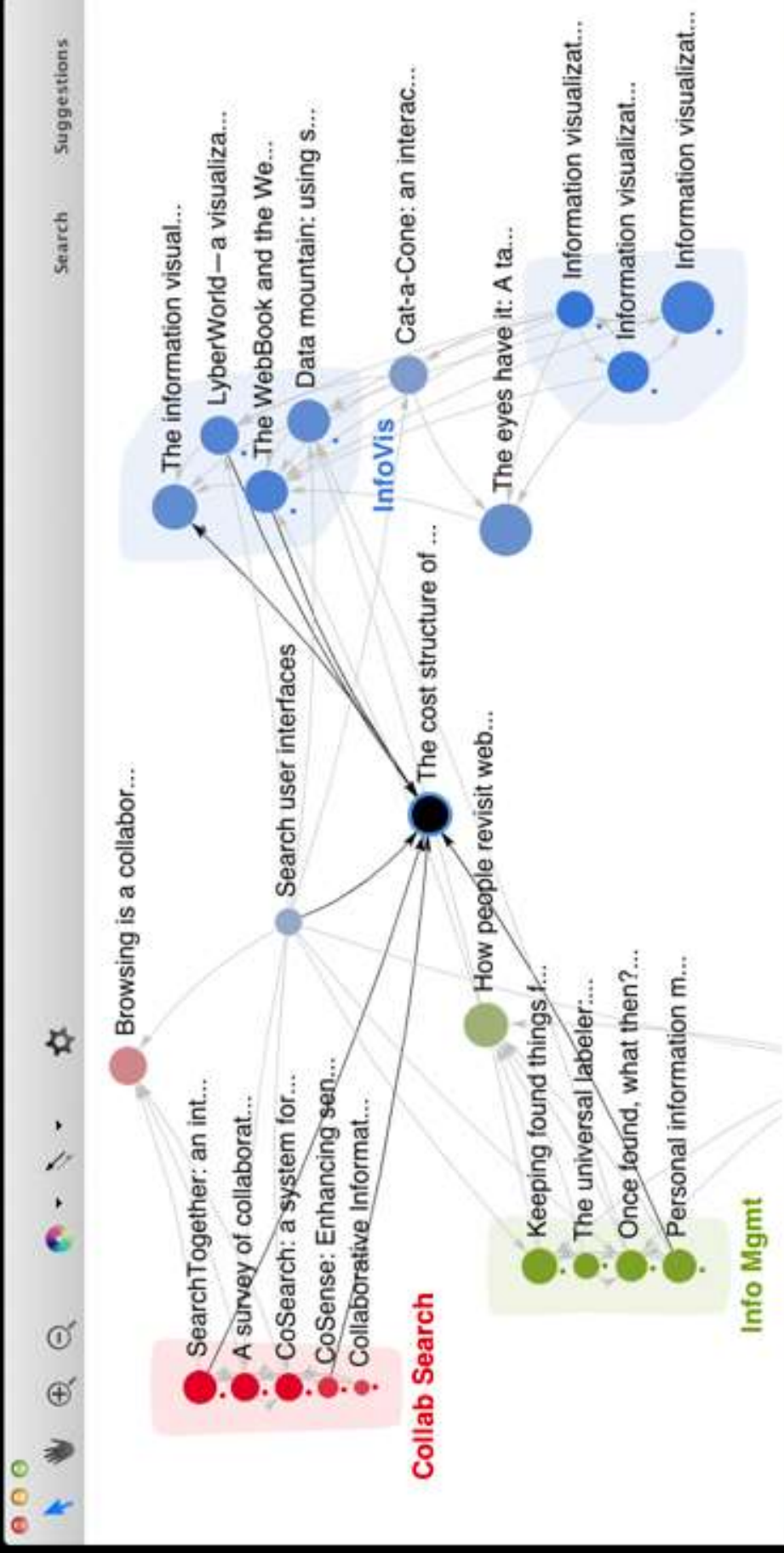


Citation network

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

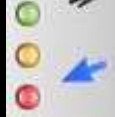


The cost structure of sensemaking

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

245 citations 8 versions

PDF 1993



Search

Suggestions

For The cost structure of sensemaking

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

The cost structure of sen...

PDF 1993

The cost structure of sensemaking

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

245 citations 8 versions

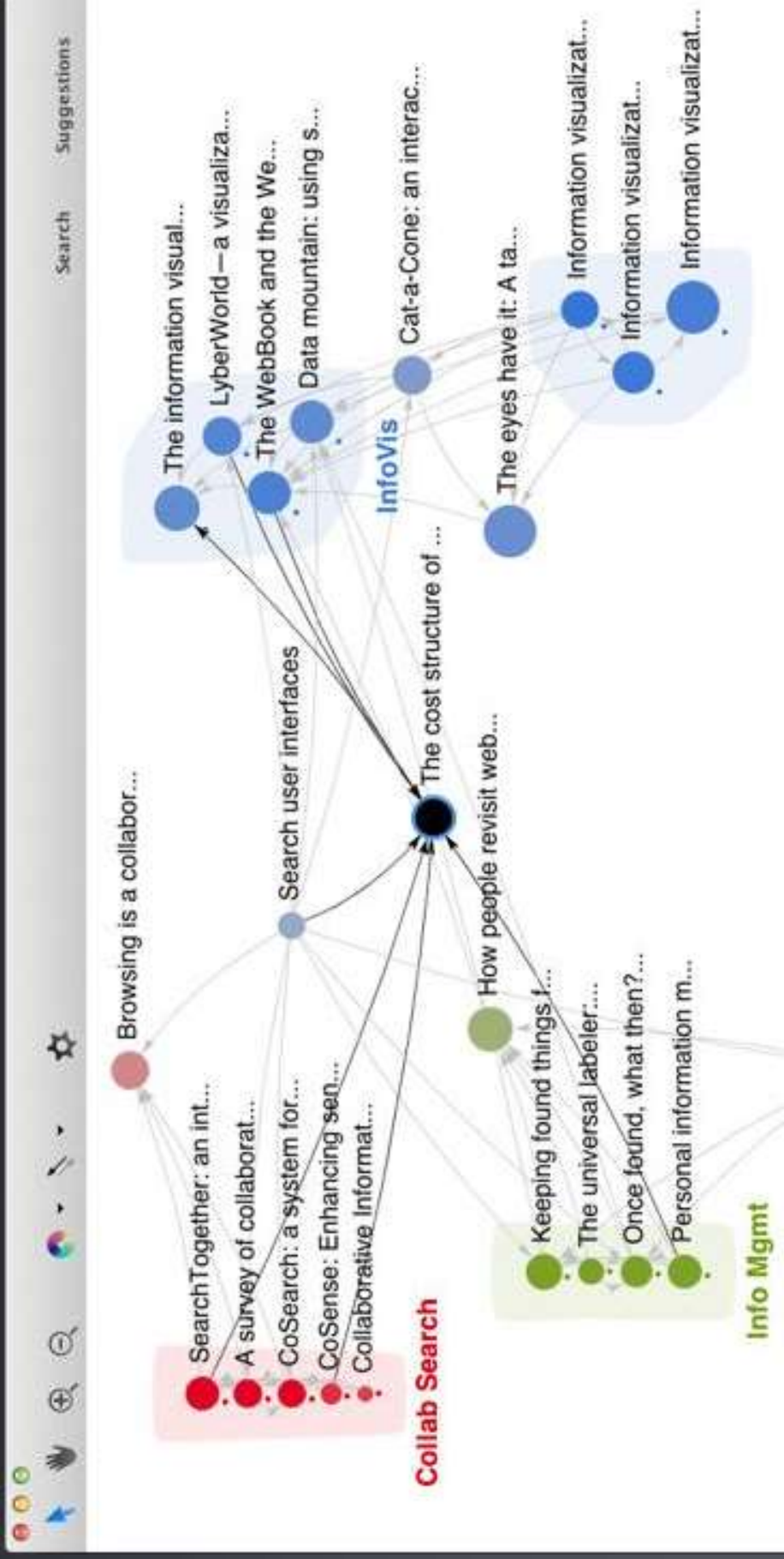


Key Ideas (Recap)



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

Cluster Data

Add Group

Recommendations:

End User Programming

End users creating effective softw...
End user software engineering: chi...
Invited research overview: end-us...
Brad A. Myers
Margaret M. Burnett
Mary Beth Rosson
Andrew Jensen Ko
Alan F. Blackwell

Show: All

Not Interested

Automatically generating user inte...
Decision-Theoretic User Interface ...
Daniel S. Weld
Krzysztof Z. Gajos
Automatically generating
Exploring the design space
Predictability and accuracy
Planning-Based Control of

Brad A. Myers
The garnet user interface developm...
Using HCI Techniques to Design a M...
Creating charts by demonstration.
The Amulet User Interface Developm...
Easily Adding Animations to Interfac...
Simplifying video editing using metad...
SILVER: simplifying video editing wit...

Text Entry

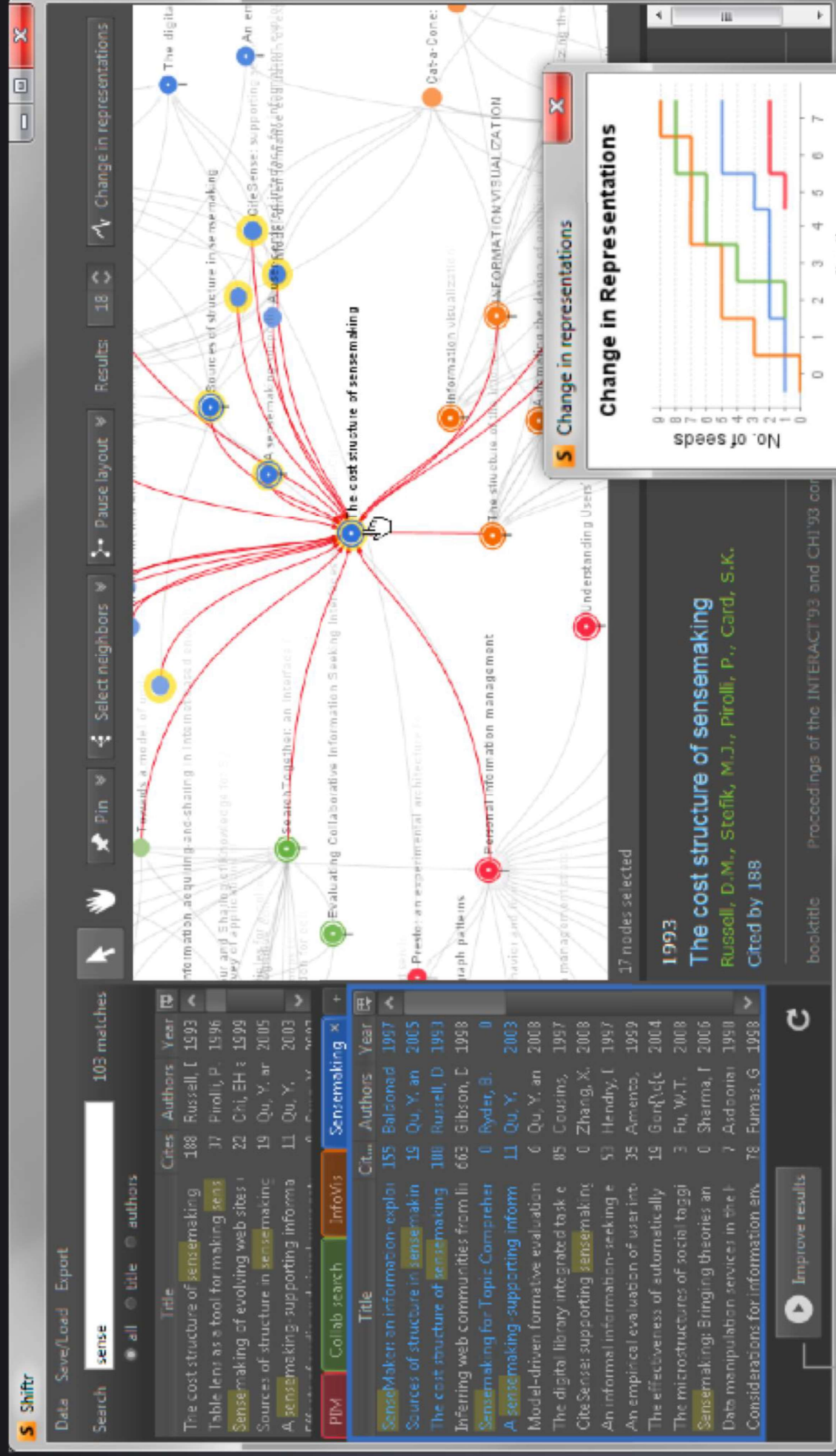
In-stroke word completion.
Integrating isometric joysticks into...
Eyes on the road, hands on the whe...
An alternative to push, press, and t...
Maximizing the guessability of symb...
Few-key text entry revisited: mnem...
Text entry from power wheelchairs: ...
Joystick text entry with date stamp, ...

Show: Papers

Interface Generation

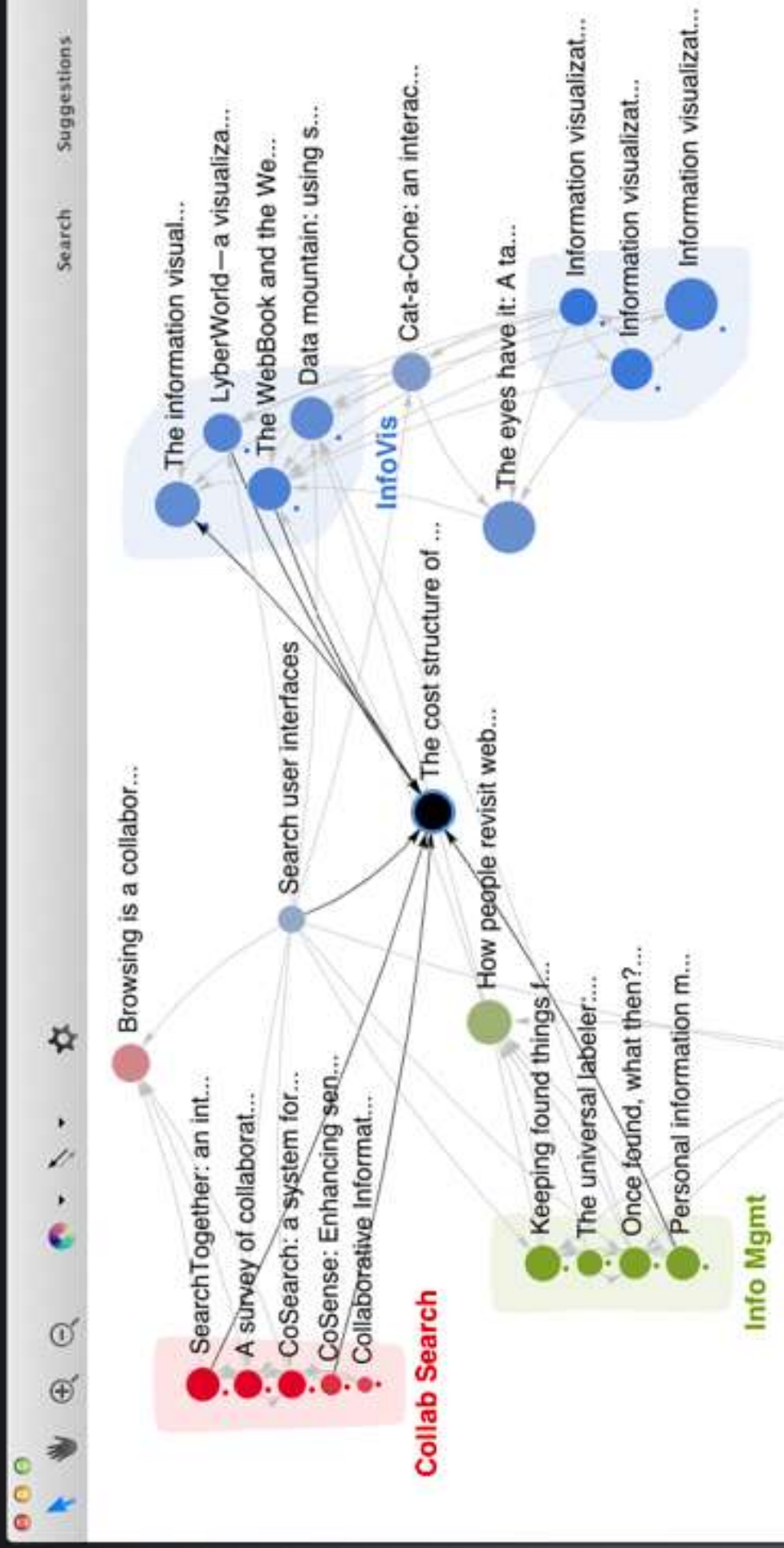
Huddle: automatically generating i...
UNIFORM: automatically generatin...
Demonstrating the viability of auto...
Jeffrey Nichols
Brandon Rothrock
Duen Horng Chau

Apollo 2010



Apollo 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 Piroli, 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

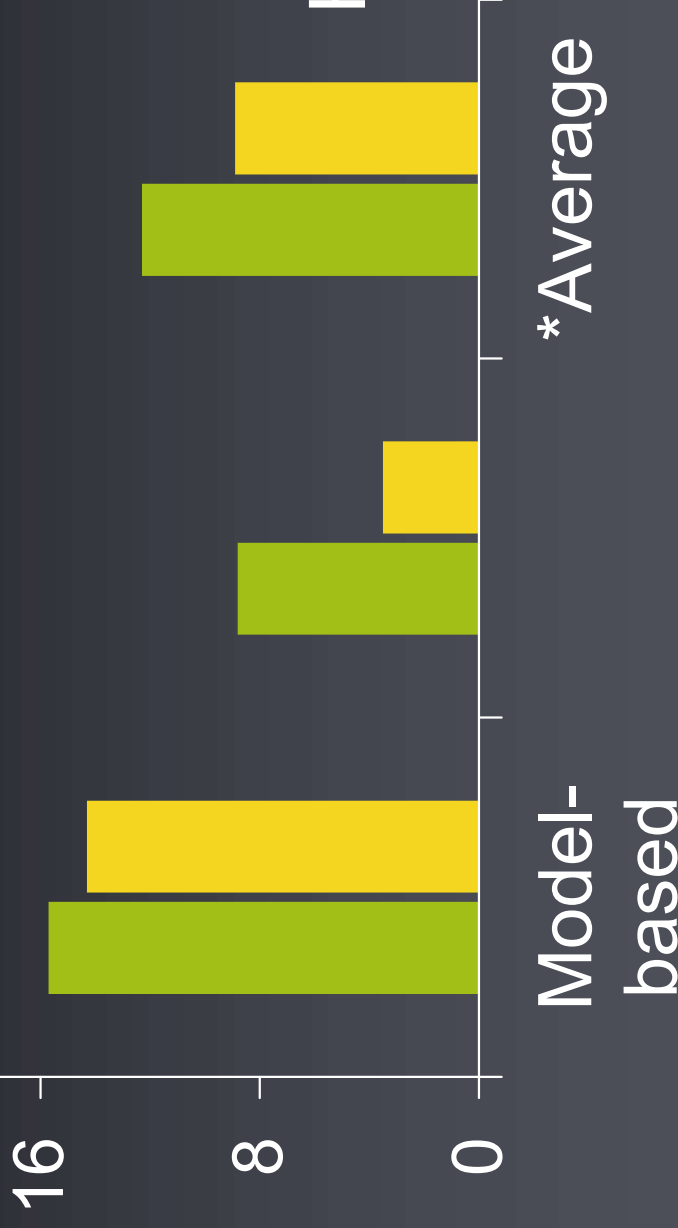
Judges'

Scores

■ Scholar

■ Apolo

Score



Higher is better.
Apolo wins.

*Average

Model-based

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

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

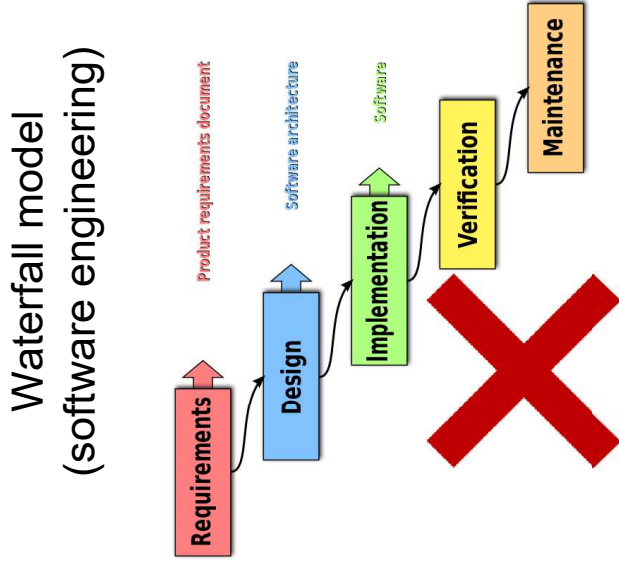
Practitioners' guide to building (interactive) applications

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**



If you want to know more about people...
<http://amzn.com/0321767535>

