

32<sup>nd</sup> Annual ACM SIGIR '09  
Boston, USA, Jul 19-23 2009



# Telling Experts from Spammers

Expertise Ranking in Folksonomies

Michael G. Noll

Christoph Meinel

Hasso Plattner Institute

(Albert) Ching-Man Au Yeung

Nicholas Gibbins

Nigel Shadbolt

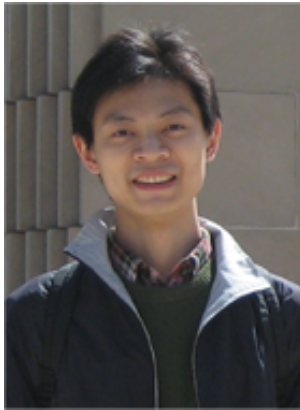
Uni Southampton

# Opening joke (under construction)



**Michael G. Noll, <http://www.michael-noll.com/>**

- Bi-national Ph.D. candidate in Computer Science at the Hasso Plattner Institute in Potsdam, Germany, and the University of Luxembourg
- Working as external doctoral student at the satellite operator SES ASTRA (Luxembourg) in the industrial R&D project “Safer Internet”
- Thesis title:  
*Understanding and Leveraging the Social Web for Information Retrieval*

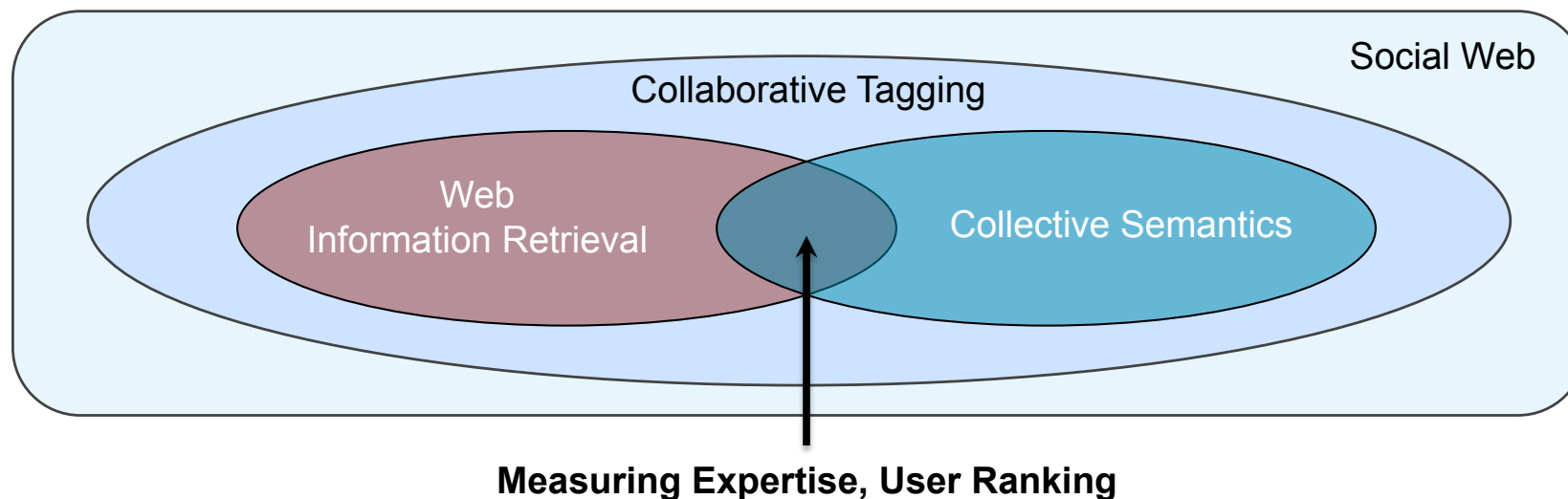


**Albert Au Yeung, <http://users.ecs.soton.ac.uk/cmayer06r/>**

- PhD candidate in Computer Science at the University of Southampton
- Previously obtained BEng (Information Engineering) and MPhil (Computer Science) from the Chinese University of Hong Kong
- Thesis title: *From User Behaviours to Collective Semantics*  
Study how implicit semantics can be harvested from social interactions on the Web, focusing on collaborative tagging as a prominent example

## Why do we work together?

5

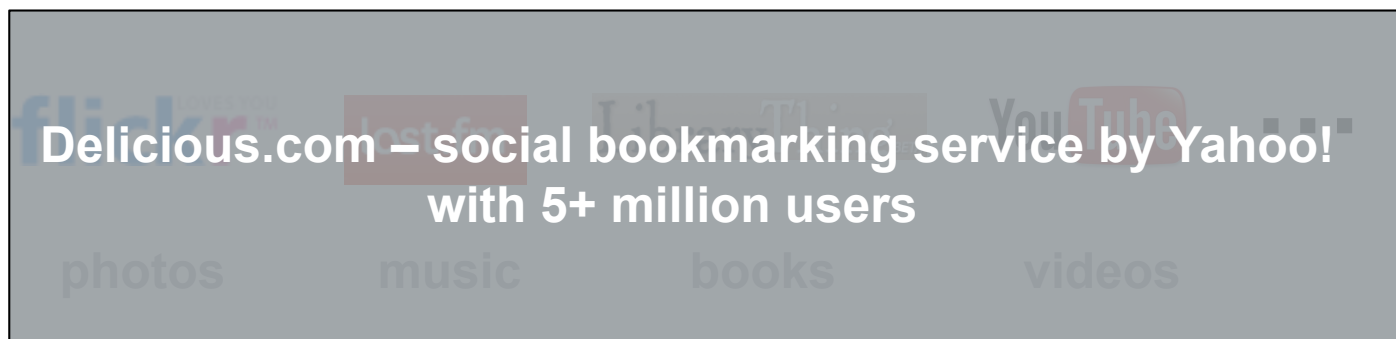


- Our common interest: ranking users according to their expertise
- In Web IR, we want to identify the experts so that we can get the best resources to satisfy our information needs
- Expertise/Trustworthiness of users is a kind of implicit quality of users that can be determined by analyzing collective user behavior
- It's FUN! 😊

# Introduction

## Folksonomies and Collaborative Tagging

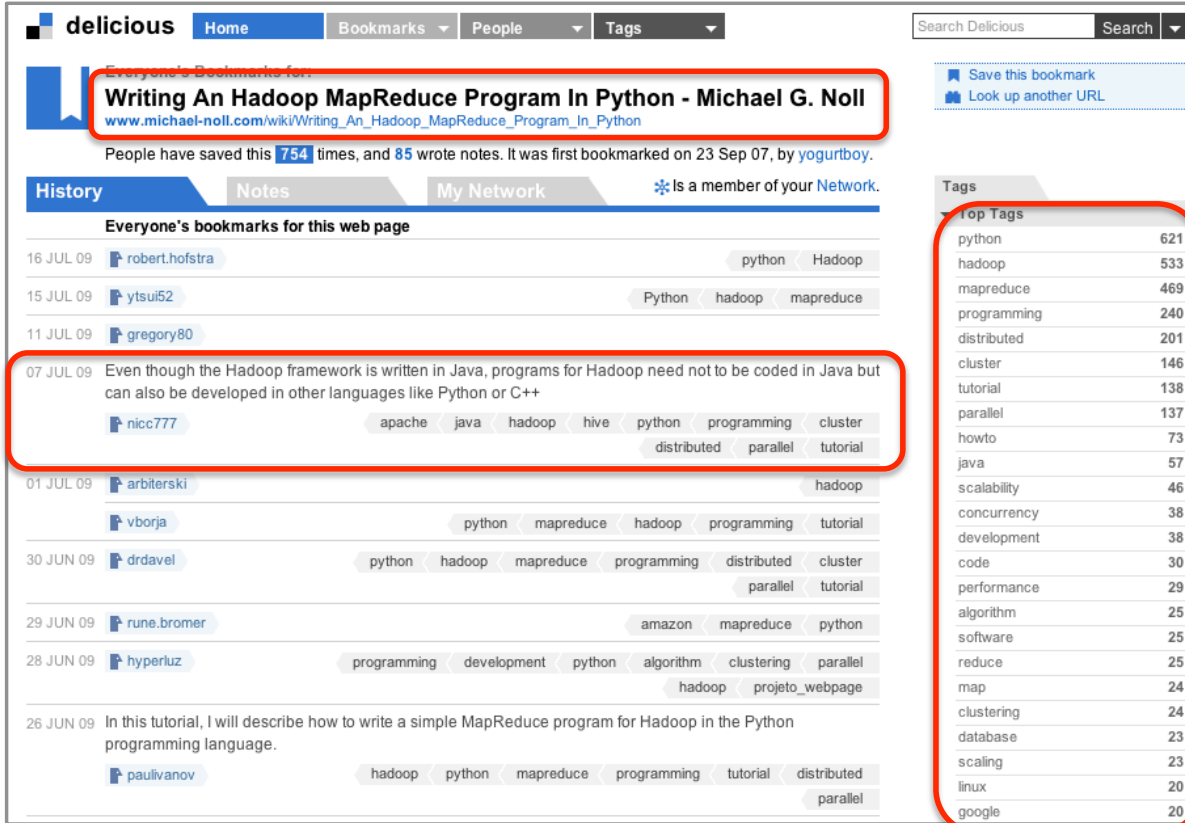
- Large and still increasing popularity in the WWW today



- Idea: Freely annotating resources with keywords aka “tags”
- Result: bottom-up “categorization” by end users, aka “folksonomy”
- Used for organizing resources, sharing, self-promotion, ...
- Additional effect: **new means** of resource/information retrieval and discovery

Web page

User bookmark



The screenshot shows a Delicious bookmark page for the article "Writing An Hadoop MapReduce Program In Python" by Michael G. Noll. The page is titled "Everyone's Bookmarks for: Writing An Hadoop MapReduce Program In Python - Michael G. Noll" and shows it has been saved 754 times. The "History" tab is active, displaying a list of users who bookmarked the page, including their names, dates, and associated tags. A red box highlights the first bookmark entry from July 7, 2009, by user "nicc777", which includes tags like "apache", "java", "hadoop", "hive", "python", "programming", "cluster", "distributed", "parallel", and "tutorial". Another red box highlights the article title and URL at the top. On the right side, a "Tags" section lists "Top Tags" with their respective counts, such as "python" (621), "hadoop" (533), and "mapreduce" (469). A red box also highlights this "Tags" section.

Tag	Count
python	621
hadoop	533
mapreduce	469
programming	240
distributed	201
cluster	146
tutorial	138
parallel	137
howto	73
java	57
scalability	46
concurrency	38
development	38
code	30
performance	29
algorithm	25
software	25
reduce	25
map	24
clustering	24
database	23
scaling	23
linux	20
google	20

Tags of all users

Example: Web page bookmarked by **754 users**, first bookmark from **09/2007**



## Two related goals for our work on expertise in folksonomies:

- 1 Identifying and promoting experts for a given topic**  
Weighting user input, giving (better) recommendations, identify trendsetters for marketing/advertising/product promotion, etc.  
***Topic := conjunction or disjunction of one or more tags***
- 2 Demoting spammers**  
Reduce impact of spam and junk input thereby improving system quality, performance, operation

# Models

## What makes an expert an expert?

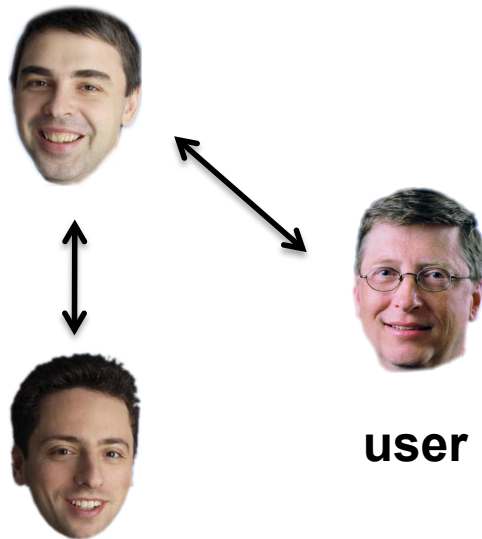
Postulation of two assumptions of expertise for resource discovery, grounded on literature from **computer science** (that's you) and **psychology**

- 1 Mutual reinforcement of user expertise and document quality**  
Expert users tend to have many high quality documents, and high quality documents are tagged by users of high expertise.
- 2 Discoverers vs. followers**  
Expert users are discoverers – they tend to be the first to bookmark and tag high quality documents, thereby bringing them to the attention of the user community. Think: researchers in academia.

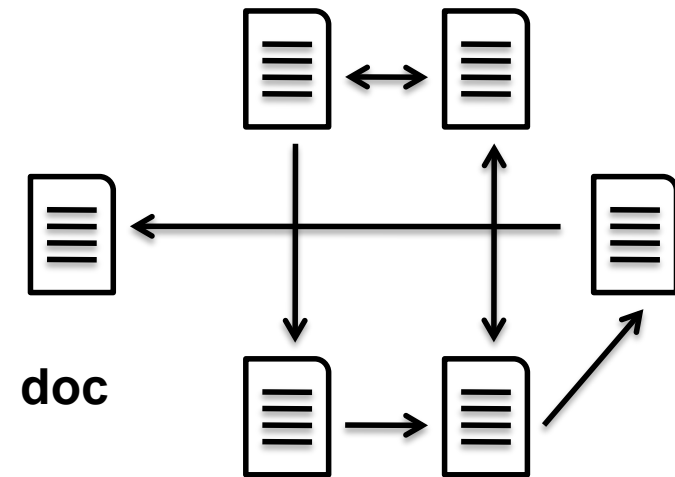
# Model of expert users

12

**user network  
(social graph)**



**document network  
(Web graph)**



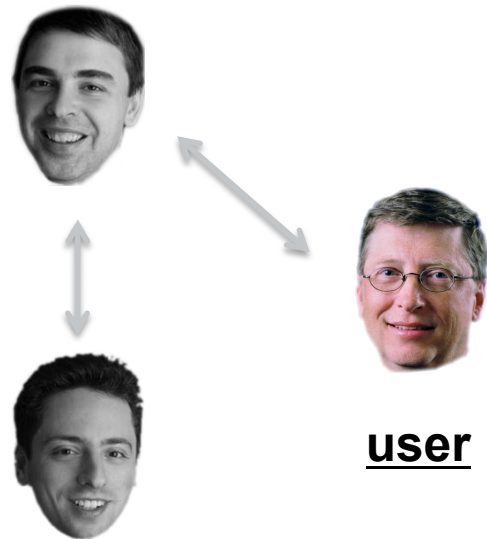
**tags**  
**timestamp**

Context of social bookmarking / collaborative tagging

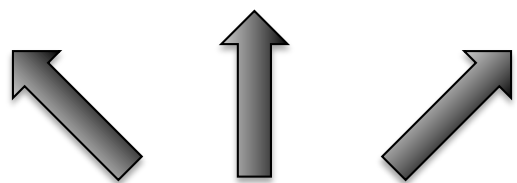
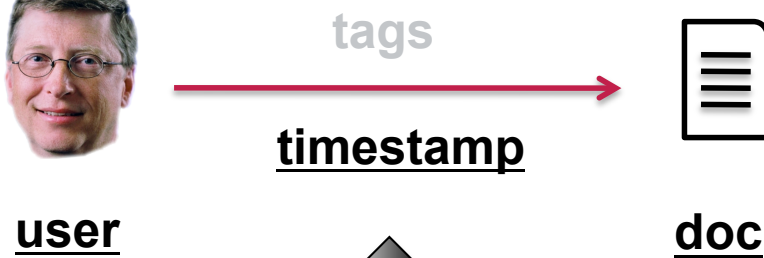
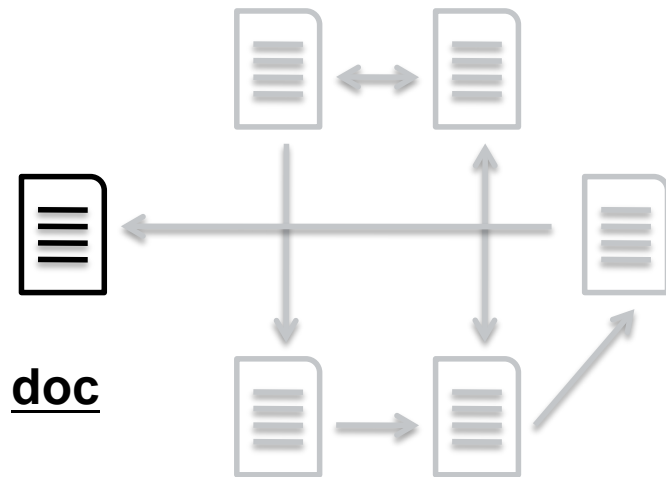
# Model of expert users

13

user network  
(social graph)



document network  
(Web graph)



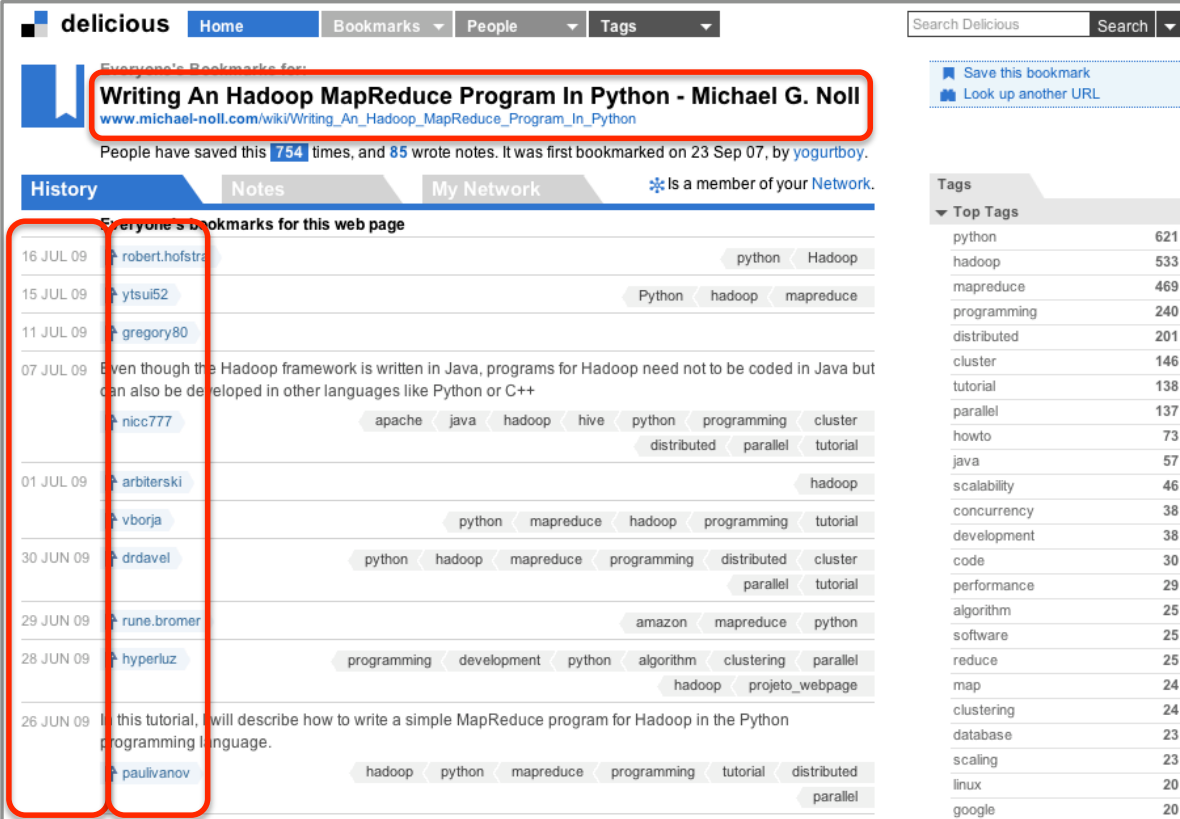
**Our Focus**

# Model of expert users

14

## Bookmarking history of a Web page

Web page



delicious Home Bookmarks People Tags

Search Delicious Search

Save this bookmark  
Look up another URL

Tags

Top Tags

python	621
hadoop	533
mapreduce	469
programming	240
distributed	201
cluster	146
tutorial	138
parallel	137
howto	73
java	57
scalability	46
concurrency	38
development	38
code	30
performance	29
algorithm	25
software	25
reduce	25
map	24
clustering	24
database	23
scaling	23
linux	20
google	20

History Notes My Network

Everyone's bookmarks for this web page

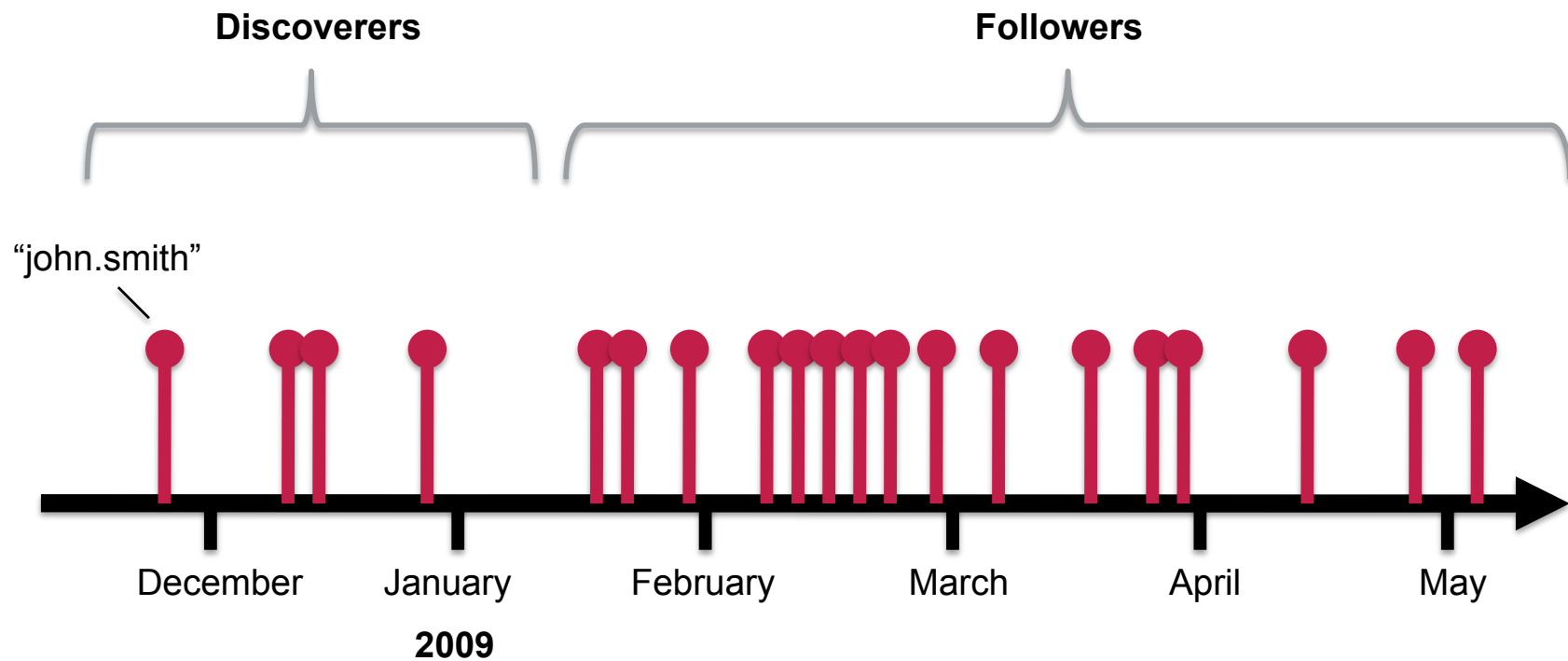
Date	User	Tags
16 JUL 09	robert.hofstra	python Hadoop
15 JUL 09	ytsui52	Python hadoop mapreduce
11 JUL 09	gregory80	
07 JUL 09	nicc777	apache java hadoop hive python programming cluster distributed parallel tutorial
01 JUL 09	arbiterski	hadoop
	vborja	python mapreduce hadoop programming tutorial
30 JUN 09	drdavel	python hadoop mapreduce programming distributed cluster parallel tutorial
29 JUN 09	rune.bromer	amazon mapreduce python
28 JUN 09	hyperluz	programming development python algorithm clustering parallel hadoop projeto_webpage
26 JUN 09	paulivanov	hadoop python mapreduce programming tutorial distributed parallel

Timeline Users

# Model of expert users

15

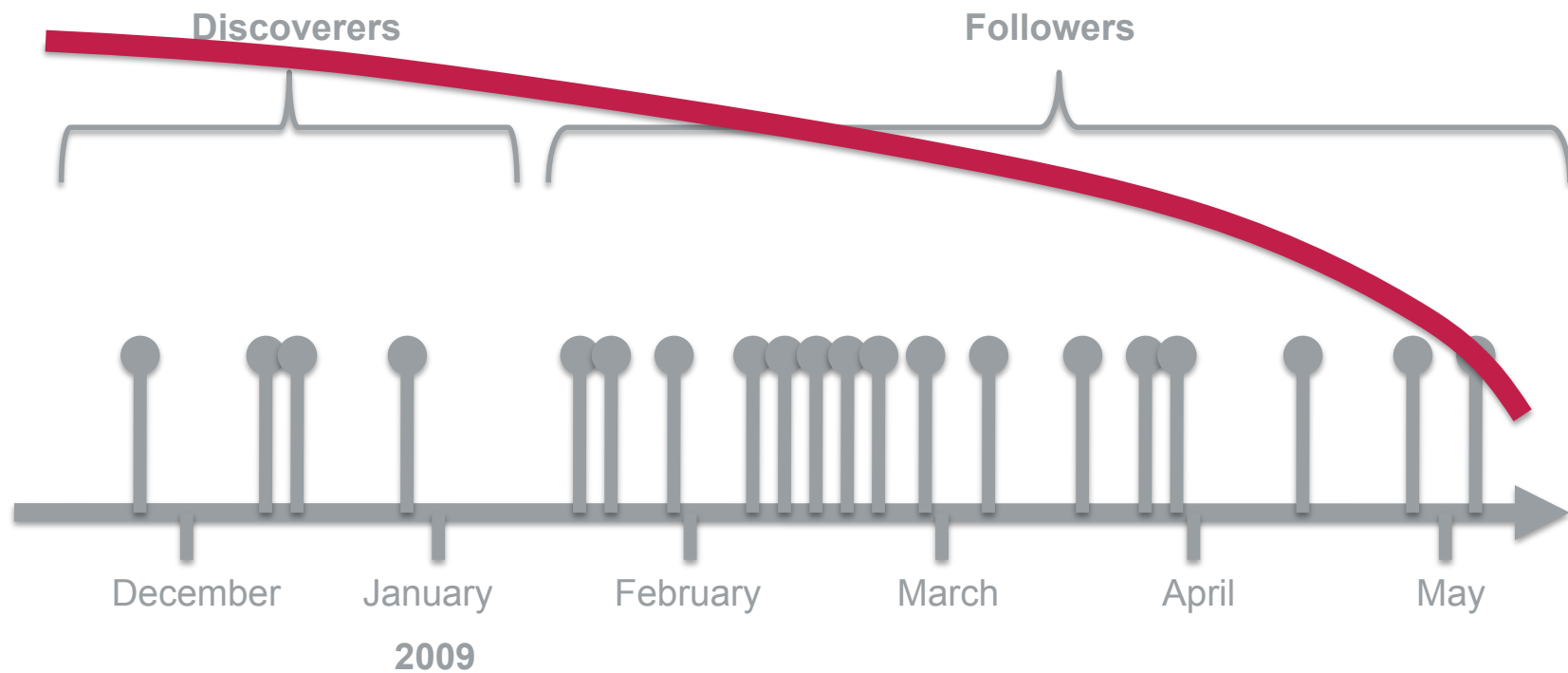
## Bookmarking history of a Web page



# Model of expert users

16

**Credit score function  $C(t)$**  → earlier discovery = more credit





# SPEAR Algorithm

## Proposed algorithm: SPEAR

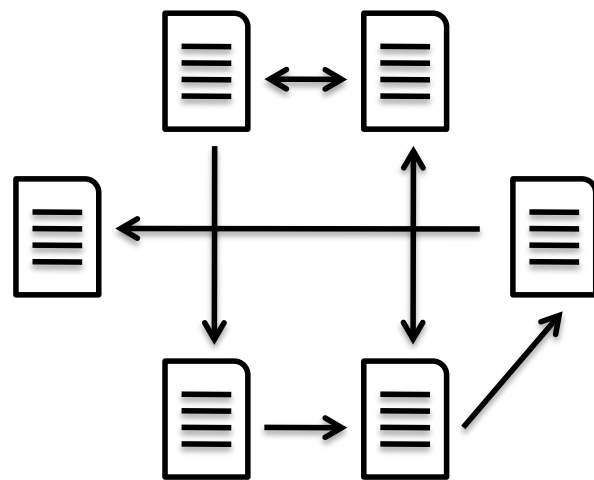
18

### **SPEAR – SPamming-resistant Expertise Analysis and Ranking**

- Based on the **HITS** (Hypertext Induced Topic Search) algorithm
  - Hubs*: pages that points to good pages
  - Authorities*: pages that are pointed to by good pages
- *Expertise and Quality* (SPEAR) similar to *Hub and Authority* (HITS)
  - Users** are **hubs** – we find useful pages through them
  - Pages** are **authorities** – provide relevant information
- Difference: only users can point (link) to pages but not vice versa

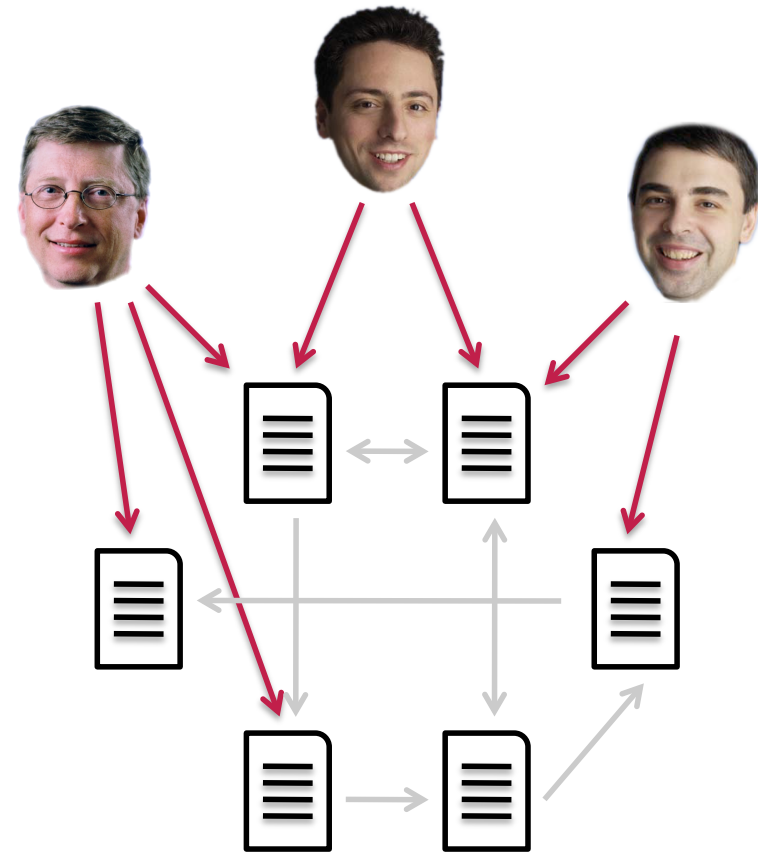
# Proposed algorithm: SPEAR

19



page ↔ page

**HITS / WWW**



user → page

**SPEAR / Folksonomy**

## Proposed algorithm: SPEAR

20

**Input**            Number of users  $M$   
                       Number of pages  $N$   
                       Set of taggings  $R_{tag} = \{ (user, page, tag, timestamp) \mid tag = tag \}$   
                       Credit score function  $C()$   
                       Number of iterations  $k$

**Output:**        Ranked list  $L$  of users by expertise in topic  $tag$

**Algorithm:**

Set  $E$  to be the vector  $(1, 1, \dots, 1) \in Q^M$

Set  $Q$  to be the vector  $(1, 1, \dots, 1) \in Q^N$

$A \leftarrow \text{Generate\_Adjacency\_Matrix}(R_{tag}, C)$

for  $i = 1$  to  $k$  do

$E \leftarrow Q \times A^T$

$Q \leftarrow E \times A$

    Normalize  $E$

    Normalize  $Q$

endfor

$L \leftarrow$  Sort users by their expertise score in  $E$

return  $L$

$E$ : expertise of users

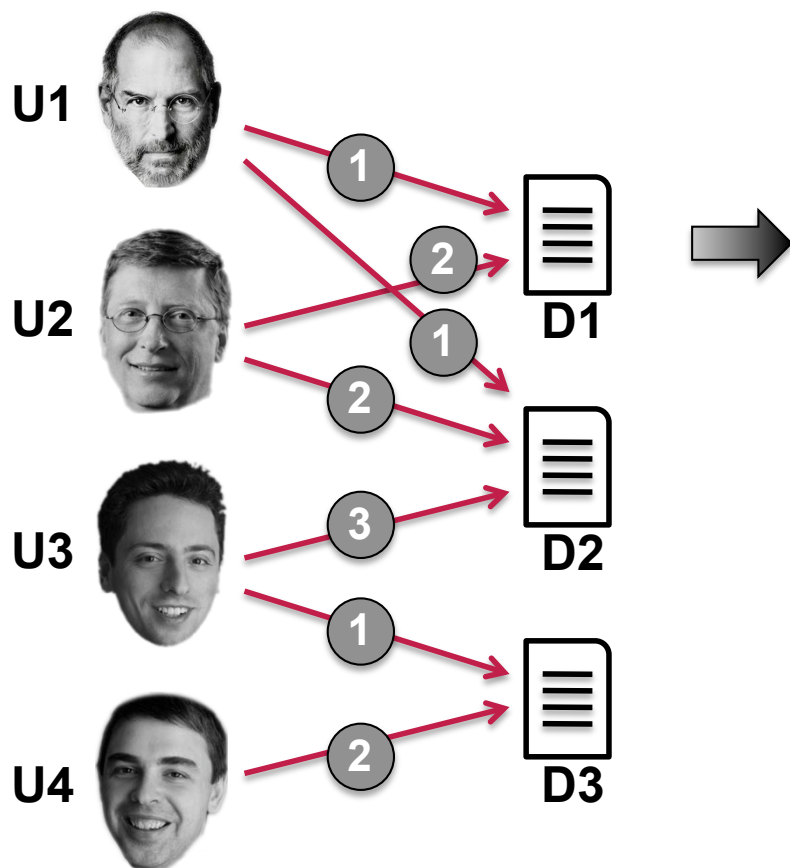
$Q$ : quality of pages

$A$ : user  $\rightarrow$  page incl. credits

mutual reinforcement  
until convergence

# Proposed algorithm: SPEAR

21



Folksonomy (simplified)

Adjacency matrix, credits applied

	D1	D2	D3
U1	1.4	1.7	0.0
U2	1.0	1.4	0.0
U3	0.0	1.0	1.4
U4	0.0	0.0	1.0



	Rank	Score
U1	1	0.422
U2	2	0.328
U3	3	0.212
U4	4	0.038

Ranked list of users by expertise

# Evaluation

## Experimental Setup

- Problem: lack of a proper ground truth for expertise
- “Who is the best researcher in this room?” 😊
- Workaround: Inserting **simulated** users into **real-world** data from Delicious.com and check where they end up after ranking
- **Real-world data set from Delicious.com** comprising 50 tags with
  - 515,000 real users (and real spammers)
  - 71,300 real Web pages
  - 2,190,000 real social bookmarks

## Experimental Setup

- **Probabilistic simulation**, simulated users generated with four parameters
  - **P1**: Number of user's bookmarks – active or inactive user?
  - **P2**: Newness – fraction of Web pages not already in data set
  - **P3**: Time preference – discoverer or follower?
  - **P4**: Document preference – high quality or low quality?



## Experimental Setup

- Simulation of 6 different user types  
Profiles (parameter values) based on recent studies + characteristics of our real-world data sets
- Experts
  - **Geek** – lots of high quality documents, discoverer
  - **Veteran** – high quality documents, discoverer
  - **Newcomer** – high quality documents, follower
- Spammers
  - **Flooder** – lots of random documents, follower
  - **Promoter** – some documents (most are his own), discoverer
  - **Trojan** – some documents, follower [next-gen spammer]

## Performance baselines

- **FREQ(UENCY)**

“Most popular” approach – simple frequency count, looks only at quantity. Seems to be the dominant algorithm in use in practice.

- **HITS**

Algorithm on which SPEAR is based. Uses mutual reinforcement but does not analyze temporal dimension of user activity.

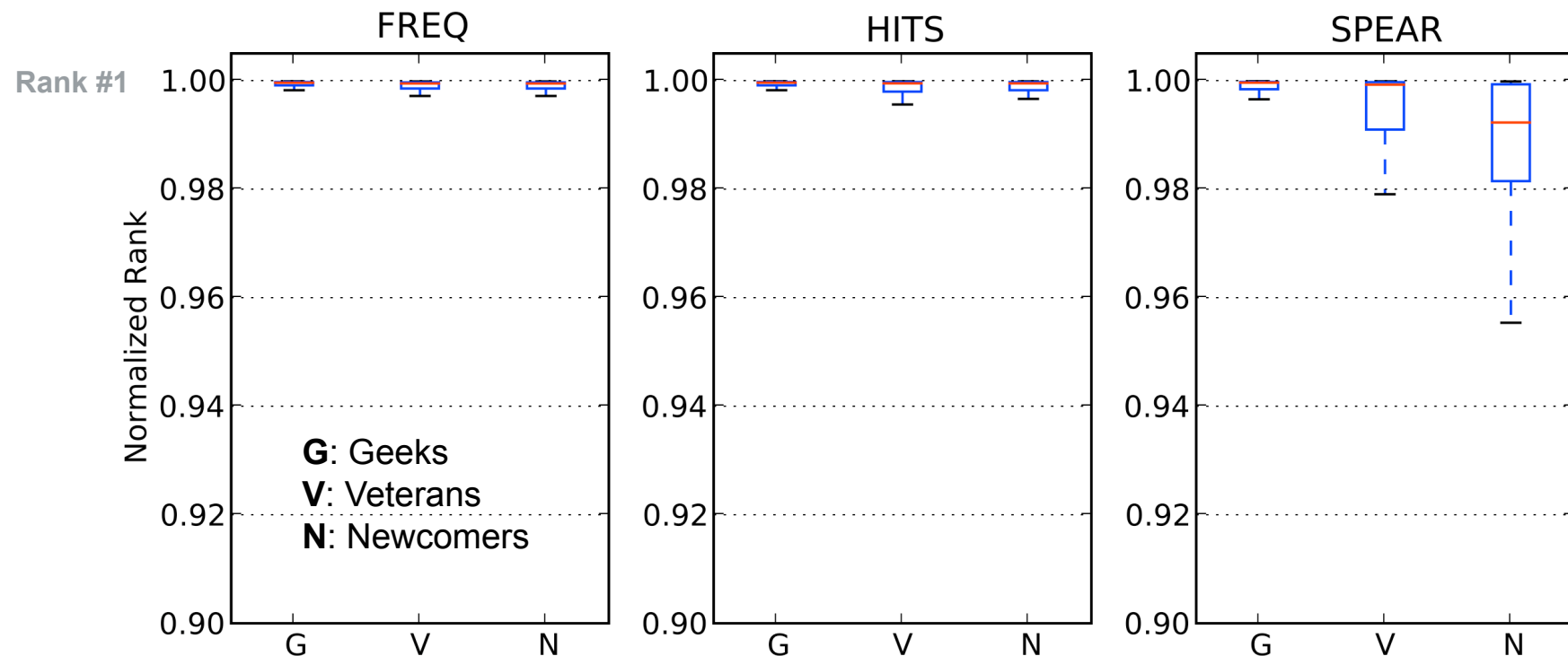
- In comparison: **SPEAR**

Uses mutual reinforcement and exploits trusted temporal data for implementing the discoverer-follower scheme.

# Experimental Results

# Experts

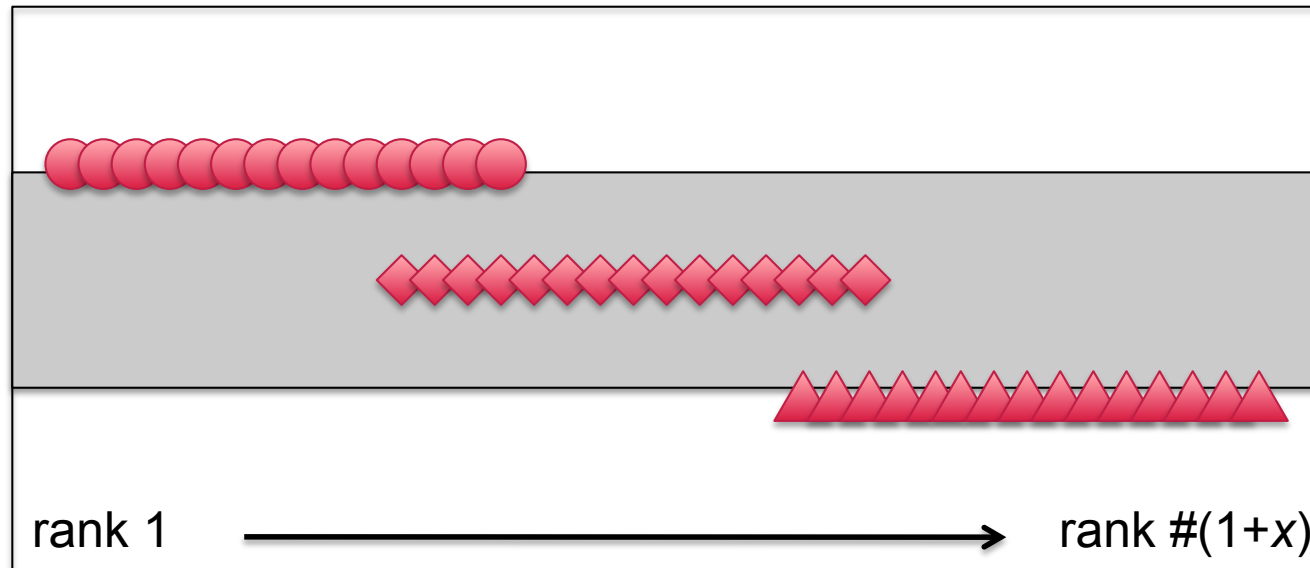
## Experts






- Only SPEAR distinguished geeks, veterans and newcomers
- FREQ and HITS clumped all expert-type users together

# A closer look

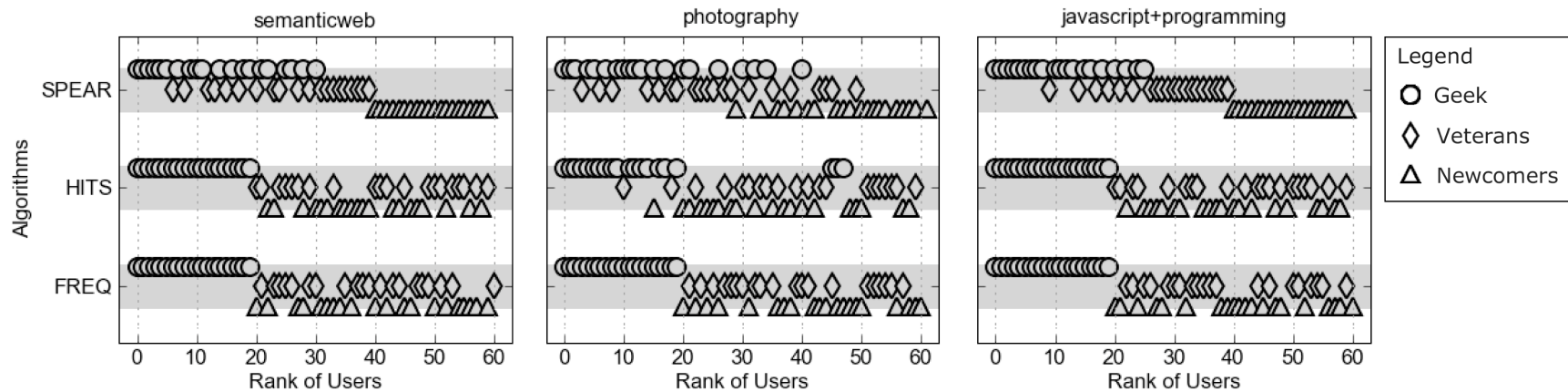
## Experts: "Ideal" result



-  Geeks
-  Veterans
-  Newcomers

**Overlaps expected due to probabilistic simulation setup**

## Experimental Results – Promoting Experts



- SPEAR differentiated all expert types better than its competitors
- SPEAR kept expected order of “geeks > veterans > newcomers”
- SPEAR was less dependent on user activity (quality before quantity)

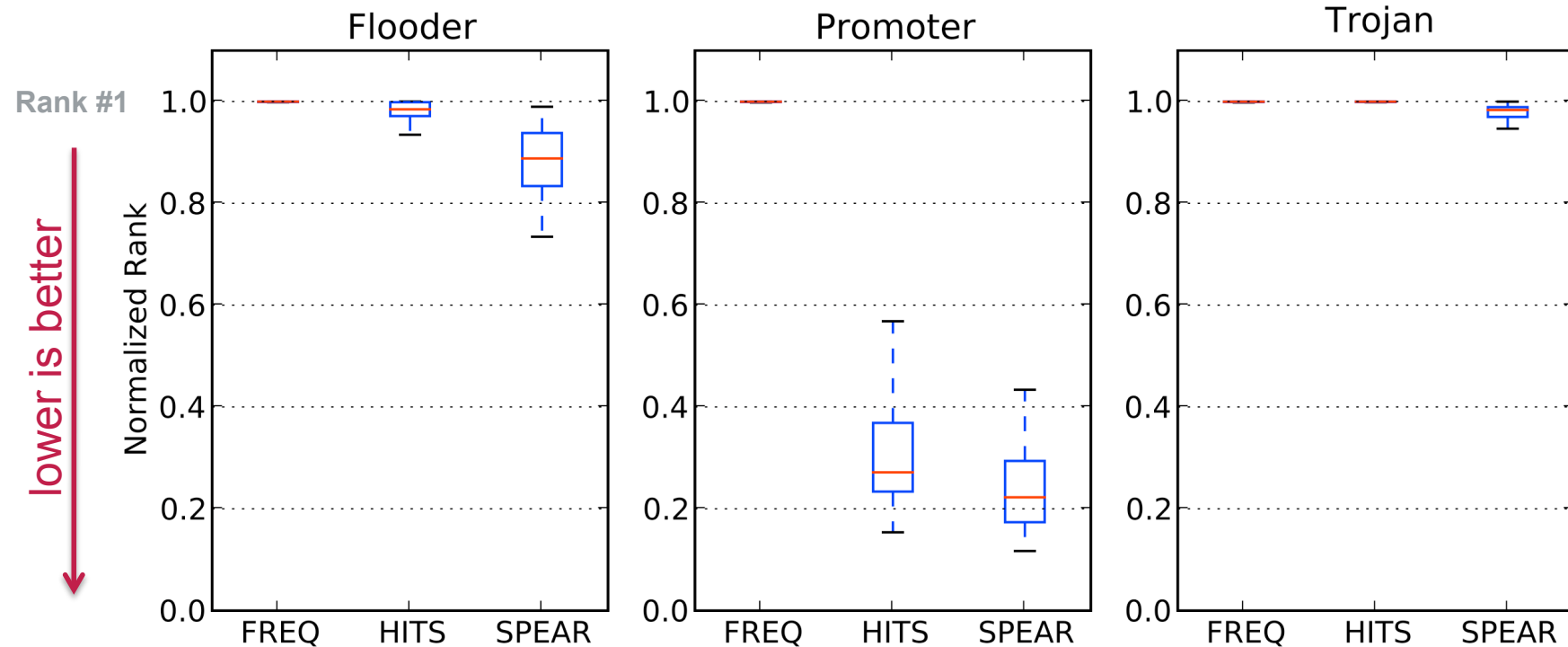


## **Qualitative analysis: manual examination of Top 10 experts for three tags “photography”, “semanticweb”, “javascript $\cap$ programming”**

- No spammers found (...phew...)
- These users seemed to be more involved or “serious” about their Delicious usage, e.g. provided optional profile information such as real name, links to their Flickr photos or microblog on Twitter
- Their number of bookmarks: from 100’s to 10,000’s
- “semanticweb”: Semantic Web researcher among the experts
- “javascript  $\cap$  programming”: Top 2 experts were professional software developers

# Spammers

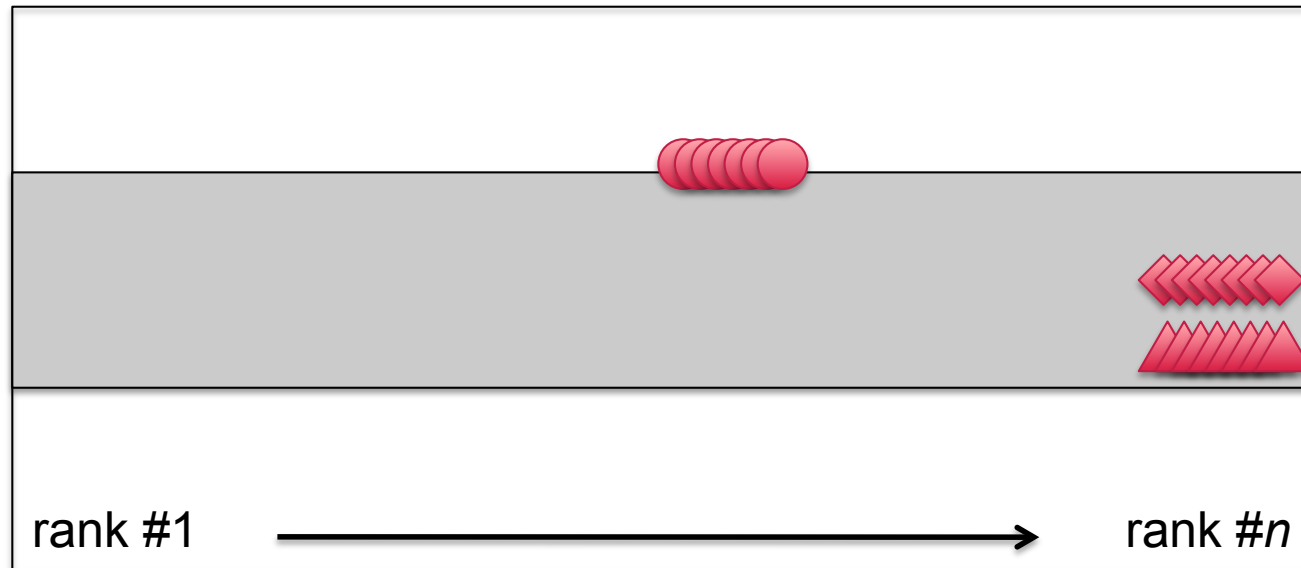
## Spammers



- SPEAR consistently outperformed FREQ and HITS
- SPEAR was the only algorithm to handle trojans (tricky spammers)

# A closer look

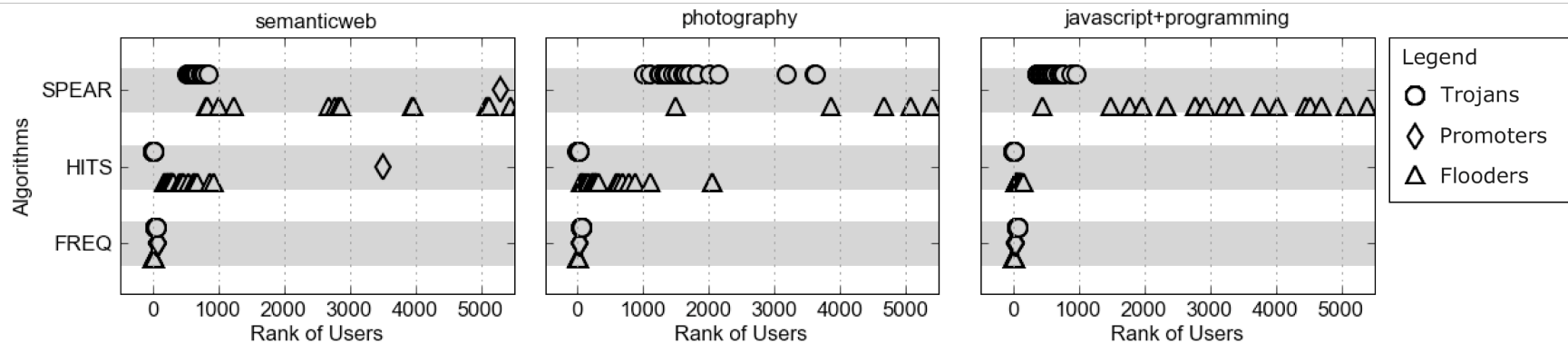
## Spammers: “Ideal” result



- Trojans
- ◆ Promoters
- ▲ Flooders

**Trojans expected to score higher because they mimic regular users for most of the time**

## Experimental Results – Demoting Spammers



- SPEAR demoted all spammer types significantly more than its competitors
- Only SPEAR demoted all trojans from the TOP 100 ranks
- FREQ completely failed to demote spammers

## **Qualitative analysis: manual examination of Top 50 users for the heavily spammed tag “mortgage” (without inserting simulated users)**

- Ranked users by their number of bookmarks = FREQ strategy
- 30 out of 50 were (real) spammers, either flooders or promoters
- Compared to FREQ, both SPEAR and HITS were able to remove these spammers from the Top 50
- SPEAR demoted spammers significantly more than HITS

## **SPEAR...**

- demoted all spammer types while still ranking experts on top
- was much less vulnerable to spammers with its reduced dependence on the activeness of the users: quality >> quantity
- *increased* difficulty for spammers to game a collaborative tagging system



# **Preliminary study: SPEAR and PageRank**

## SPEAR and PageRank

- Second SPEAR outcome: **document quality score**
- Relationship to other document quality / popularity measures?

## Questions

- “Correlation between SPEAR (folksonomy) and PageRank (Web graph)?”
- “Are documents ranked high by SPEAR also ranked high by PageRank?”

## SPEAR and PageRank

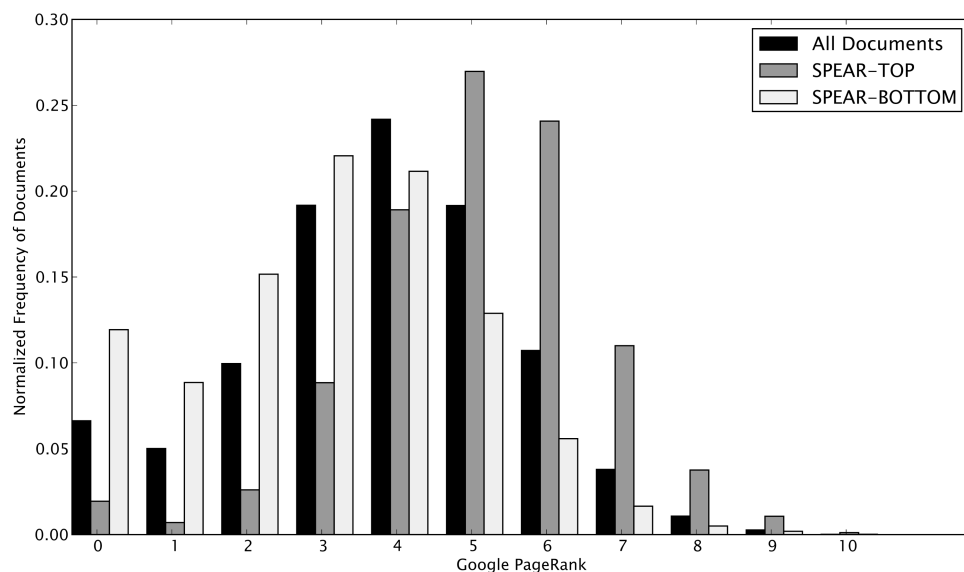
### Experiments

- We built three data sets and compared PageRank distributions:
  - *ALL* = all documents from our **n** real-world data sets
  - *SPEAR-TOP* = joint set of SPEAR Top 100 docs of all **n** data sets
  - *SPEAR-BOTTOM* = joint set of SPEAR Bottom 100 docs ...

# SPEAR and PageRank

44

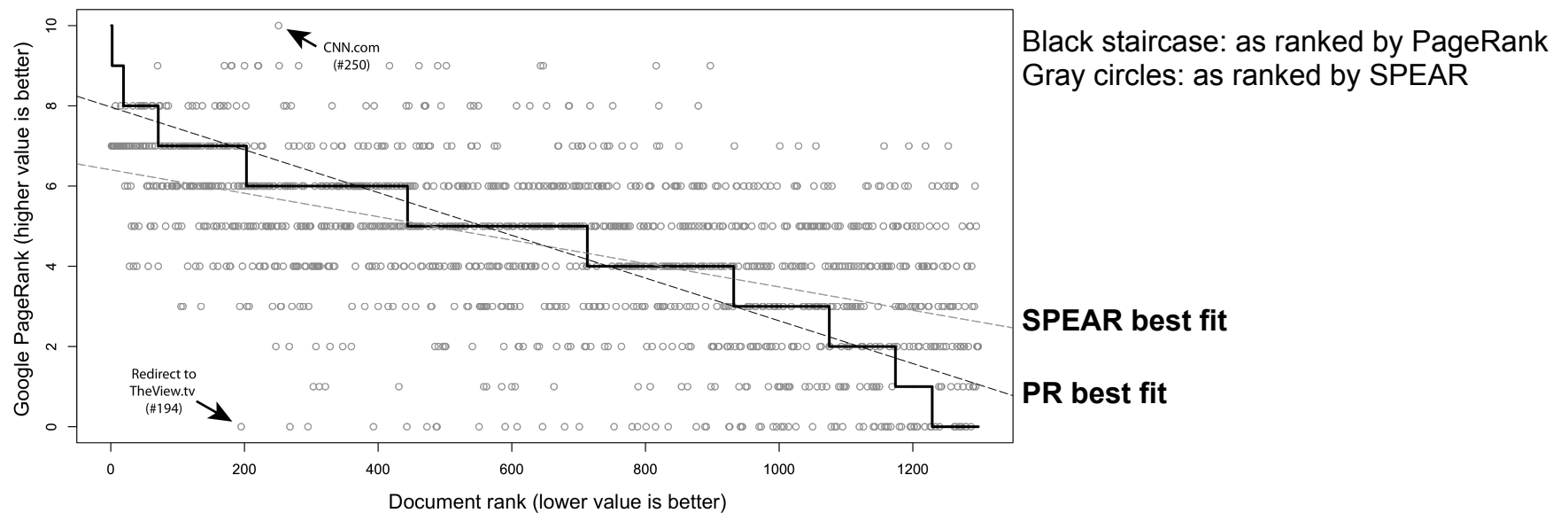
## SPEAR and PageRank



Documents	Mean PR	std.dev.	Median PR
ALL	3.71	1.81	4
SPEAR-TOP	5.05	1.61	5
SPEAR-BOTTOM	3.05	1.81	3

- Documents ranked higher (lower) by SPEAR tend to have higher (lower) PageRanks
- Mean Pearson- $r$  correlation coefficient averaged over **all** data sets:  $r = +0.324$

## PR distributions of exemplary data set “entertainment”



- Still, SPEAR generally behaves quite different from PageRank!
  - Best PR0 document [SPEAR #194] > best PR10 document [SPEAR #250]:  
PR0 document redirects to a PR8 document (homepage of TV show “The View”)

# Summary

## Conclusions

- Described a model of expertise in folksonomies for resource discovery
- Proposed an expertise ranking algorithm that is resistant to spammers
- Demonstrated how simulation techniques can be used for evaluation

## Future Work

- Quality score of Web pages deserve more investigation
- Transfer to new problem domains, e.g. blogosphere or music
- Follow-up with user & item recommendation, trend detection



Michael G. Noll  
michael.noll@hpi.uni-potsdam.de  
Hasso Plattner Institute, LIASIT



Albert Au Yeung  
cmay06r@ecs.soton.ac.uk  
University of Southampton

# Backup Slides



## Resource retrieval ~ information retrieval

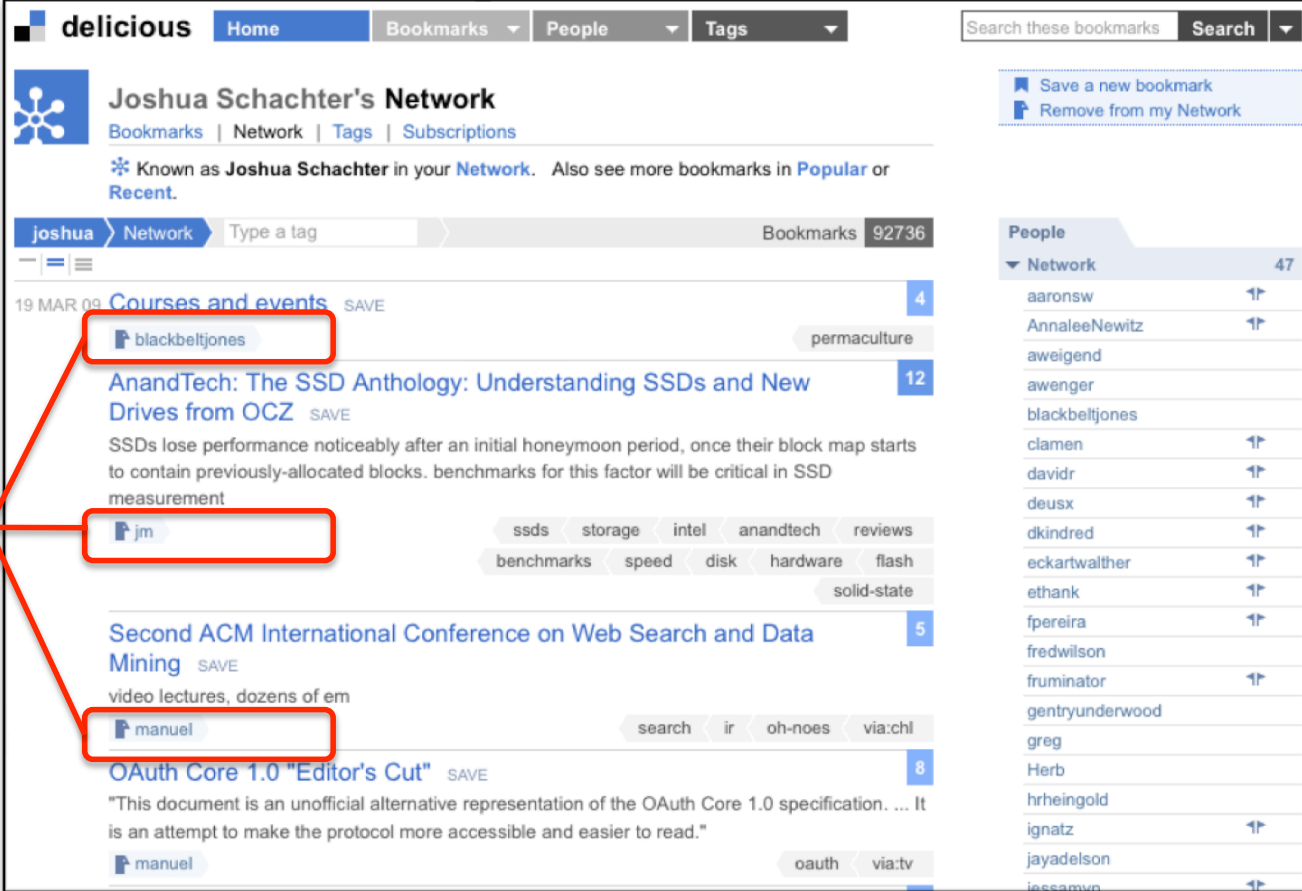
- Two types of resource discovery in collaborative tagging:
  1. Following the **tags**: subscribe or search tags to find relevant resources
  2. Following the **users**: subscribe to user feeds and receive notifications
  
- Following **expert users** provides more benefits
  - Should know the best resources with respect to a given topic
  - Should be quick in discovering and identifying new resources

# Motivation

50

Example: The user network of Joshua Schachter, founder of Delicious.com

Joshua “follows” these users and their activity



The screenshot shows the 'delicious' website interface for 'Joshua Schachter's Network'. The page displays a list of bookmarks with user avatars and names. Three users are highlighted with red boxes: blackbeltjones, jm, and manuel. A red arrow points from the text 'Joshua follows these users and their activity' to these three users.

Network	Count
aaronsw	4
AnnaleeNewitz	1
aweigend	1
awenger	1
blackbeltjones	1
clamen	1
davidr	1
deusx	1
dkindred	1
eckartwalther	1
ethank	1
fpereira	1
fredwilson	1
fruminator	1
gentryunderwood	1
greg	1
Herb	1
hrheingold	1
ignatz	1
jayadelson	1
jessamy	1