

Scalable Collaborative Filtering for Mining Social Networks

Edward Chang

Google Research, Beijing

<http://infolab.stanford.edu/~echang/>

December 12th, 08

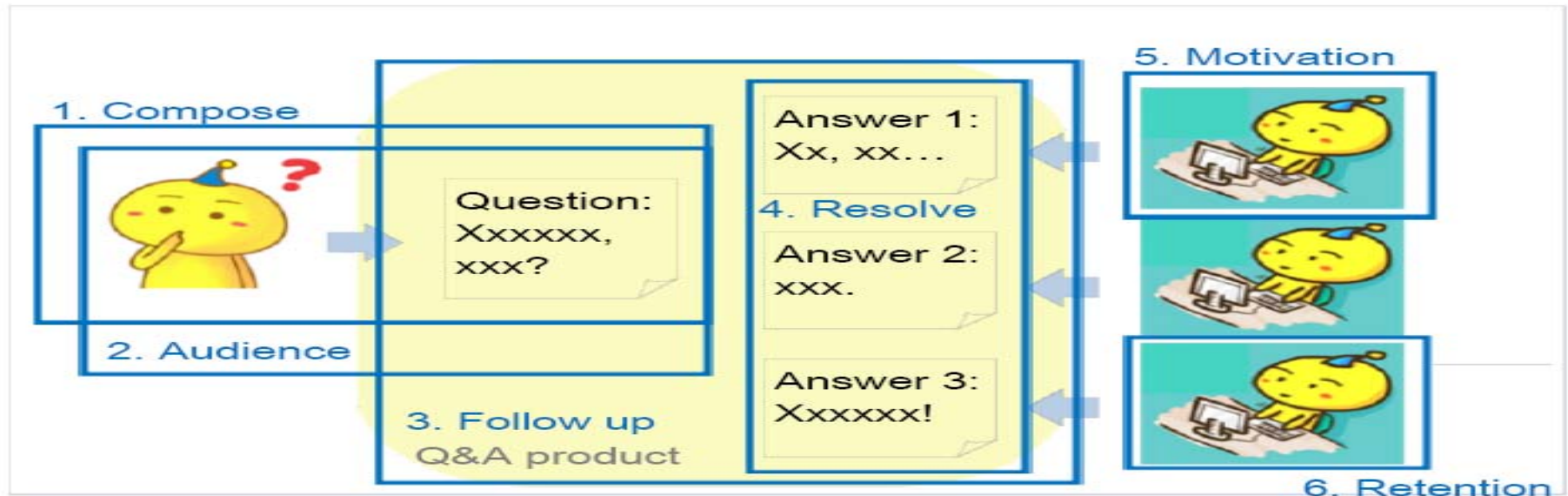
NIPS Beyond Search



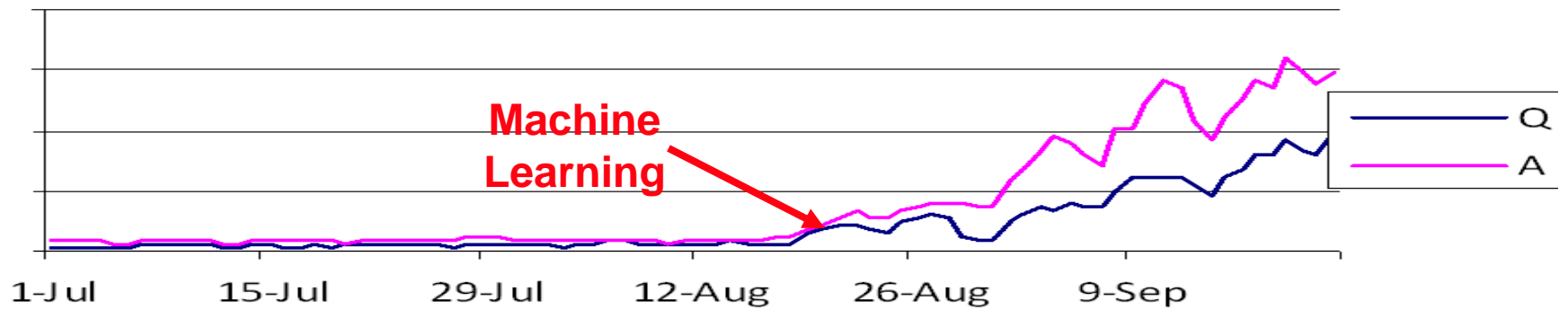
Collaborators

- Prof. Chih-Jen Lin (NTU)
- Hongjie Bai (Google)
- Wen-Yen Chen (UCSB)
- Jon Chu (MIT)
- Haoyuan Li (PKU)
- Yangqiu Song (Tsinghua)
- Matt Stanton (CMU)
- Yi Wang (Google)
- Dong Zhang (Google)
- Kaihua Zhu (Google)
- Confucius Team led by Jim Deng (Google Beijing)
- OpenSocial Team led by David Glazer (Google MTV)

Confucius, a Q&A System



Confucius Growth

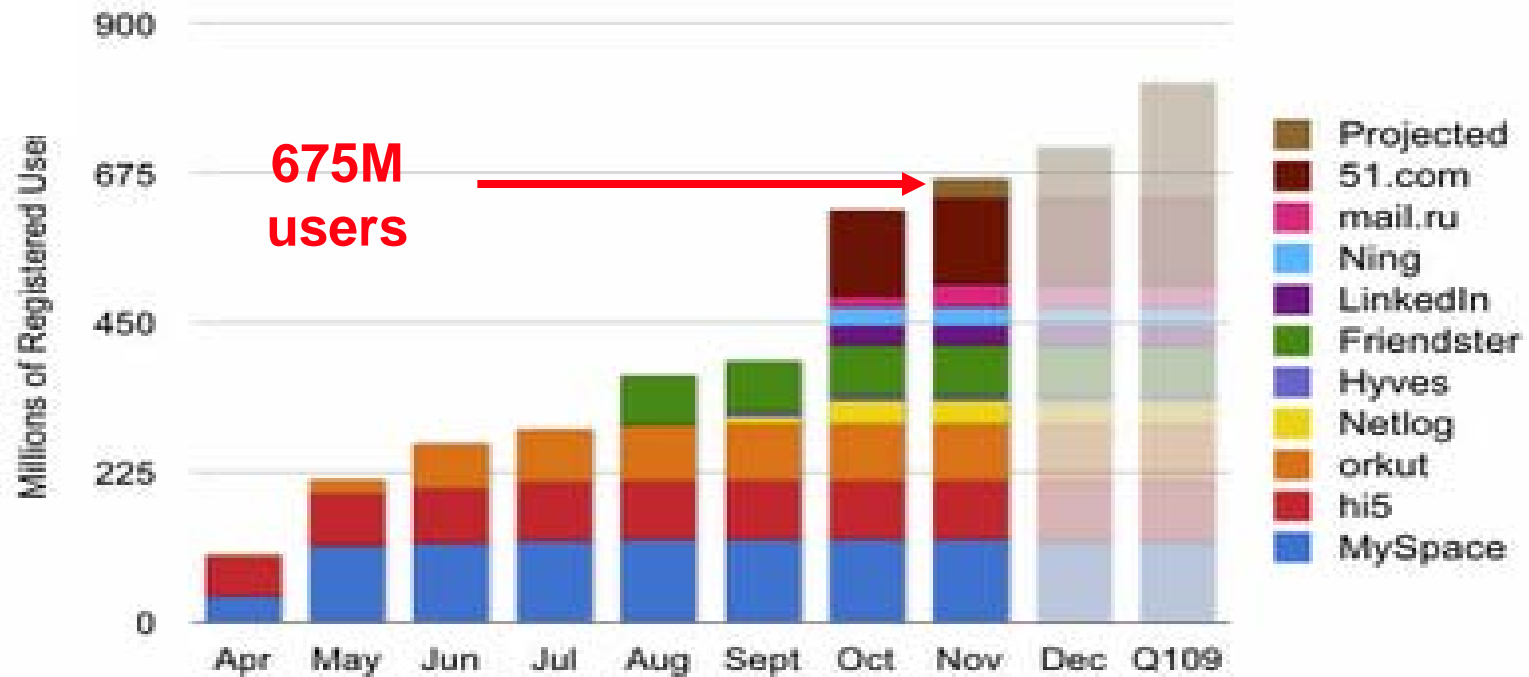


December 12th, 08

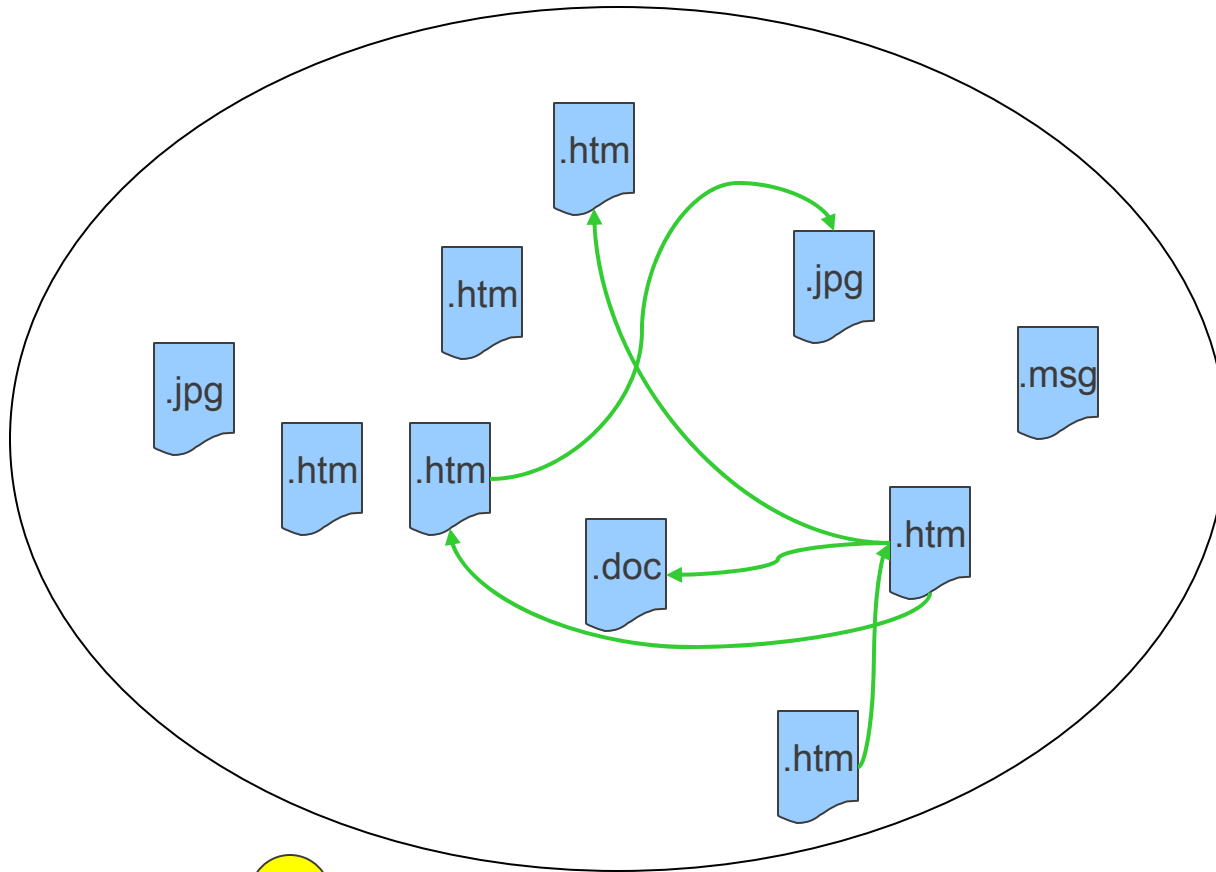
NIPS Beyond Search

3

OpenSocial



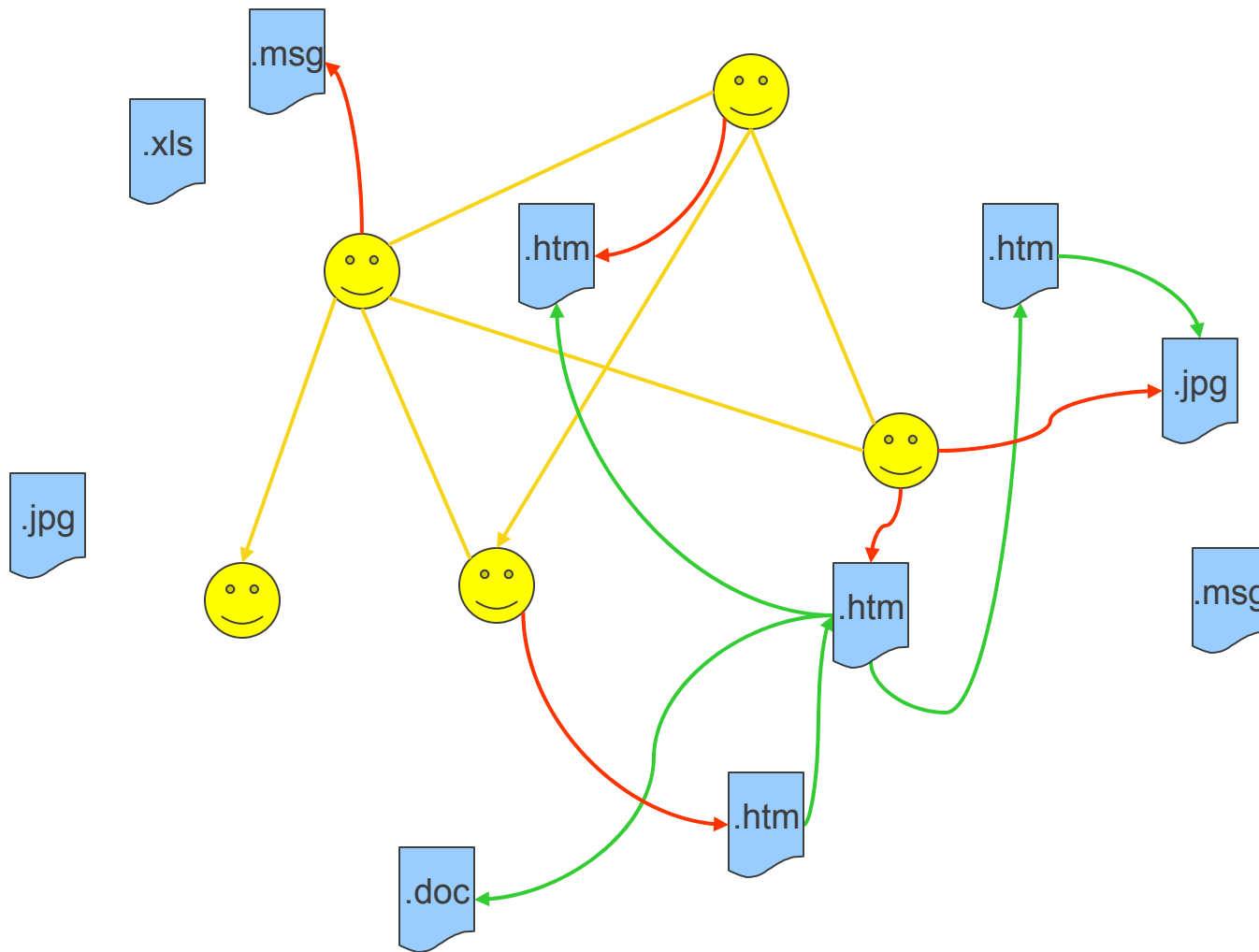
Web 1.0



December 12th, 08

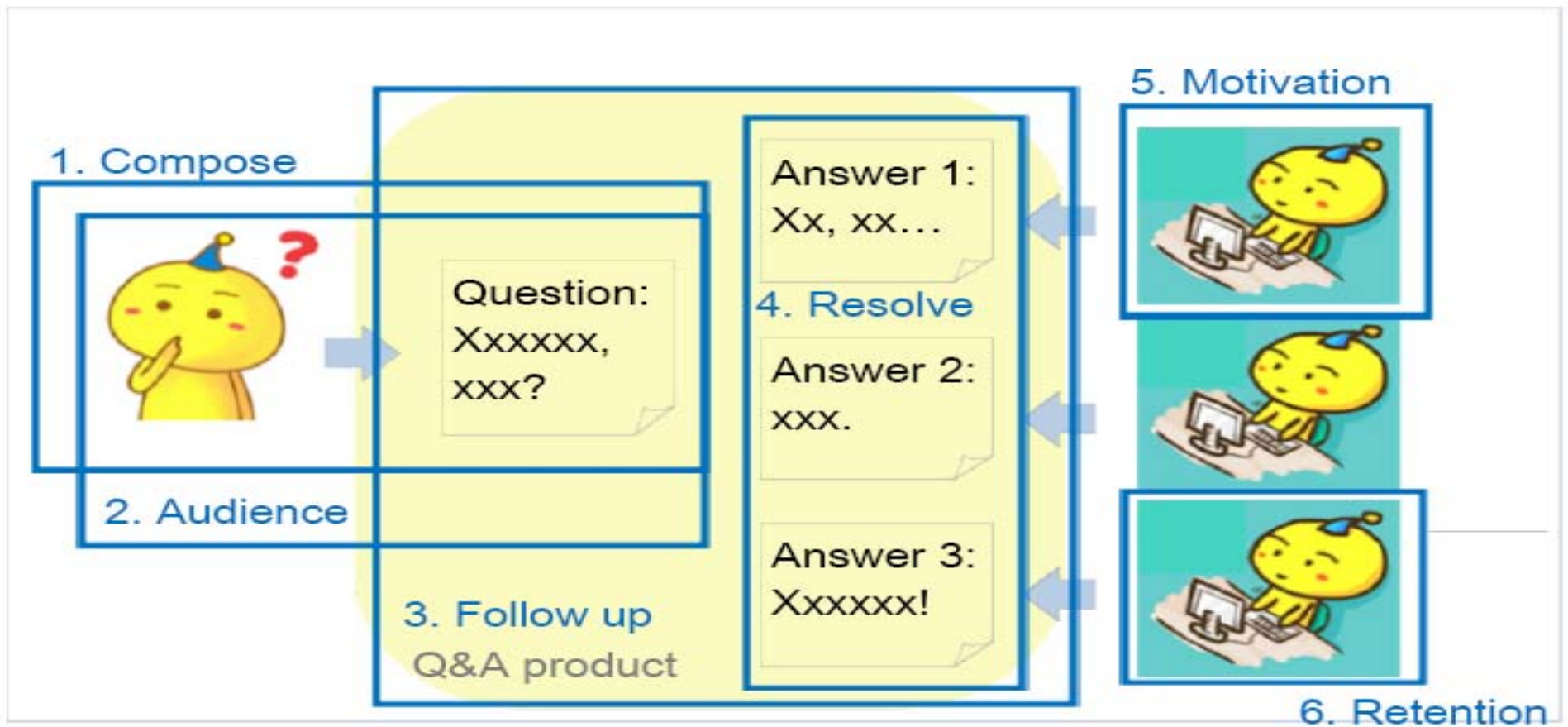
NIPS Beyond Search

Web 2.0 --- Web with People



Confucius, a Q&A system

- Allowing people to ask questions for information that cannot be found by Web search



Query: *What are must-see attractions at Yellowstone*

Google Search [Advanced Search](#) [Preferences](#)

Web Results 1 - 10 of about 12,000 for **What are must-see attractions at Yellowstone**. (0.18 seconds)

[Three Must See Attractions at Yellowstone National Park « The View ...](#)

Jan 15, 2008 ... Smith presents **Three Must See Attractions at Yellowstone** National Park posted at The View West. Interested in **Yellowstone** National Park? ... [theviewwest.com/2008/01/15/three-must-see-attractions-at-yellowstone-national-park/](#) - 26k - [Cached](#) - [Similar pages](#)

[Three Must See Attractions At Yellowstone National Park](#)

Jan 15, 2008 ... **Three Must See Attractions At Yellowstone** National Park. [ezinearticles.com/?Three-Must-See-Attractions-At-Yellowstone-National-Park&id=929265](#) - 47k - [Cached](#) - [Similar pages](#)

[Yellowstone National Park: Top Ten Attractions](#)

YELLOWSTONE NATIONAL PARK by **Yellowstone** Net. Top 10 Things to See in YNP What are the "**Must See**" attractions to view in **Yellowstone**? Start here! ... [www.yellowstone.net/topten.htm](#) - 16k - [Cached](#) - [Similar pages](#)

[Yellowstone Must-see Attractions](#)

Yellowstone's Must-See Attractions. The locations of all sites listed below are shown on the map that you receive as you enter the park. ... [www.geocities.com/dmonteit/must_see.html](#) - 8k - [Cached](#) - [Similar pages](#)

[What to See in Yellowstone](#)

Must-See Attractions -- Text Only Version · Upper Geyser Basin and Old Faithful · Grand Canyon of the **Yellowstone** · Fountain Paint Pots Trail · Wildlife ... [www.geocities.com/dmonteit/whattosee.html](#) - 10k - [Cached](#) - [Similar pages](#)
[More results from www.geocities.com »](#)

[Must See in Yellowstone National Park](#)

Query: *What are must-see attractions at Yellowstone*



At first glance, Mammoth Hot Springs appear as a frozen waterfall. Large terraces abound while being connected by trickling water. The hot acidic water from the thermal aspect below ascends through ancient limestone deposits in the area. As the water dissolves the limestone, it is carried to the surface. When the suspension cools and becomes less acidic at the surface it forms the pools and the cascading features. This area is truly an amazing and dynamic work of art.

Wildlife



- o [The Church of Jesus Christ of Latter Day Saints](#)
- o [The View West Bookstore](#)
- o [WordPress.com](#)
- o [WordPress.org](#)

ARCHIVES

- o [May 2008 \(1\)](#)
- o [March 2008 \(1\)](#)
- o [February 2008 \(15\)](#)
- o [January 2008 \(19\)](#)

BLOG STATS

o 4,702 hits

TAGS

Avalanche

[avalanche deaths](#)

[avalanche fatalities](#)

[baseball Bill Richardson bonneville dam Book Reviews](#)

[California budget](#)

[California Deficit education cuts](#)

[Election 2008 full day](#)

[kindergarten geysers goose](#)

[gossage gossage governor](#)

[Schwarzenegger hall of fame](#)

[highway 66 idaho snow jaycee](#)

[carroll kindergarten lava dome](#)

[LDS church montana](#)

[avalanche Mount St.](#)

Query: *What are must-see attractions at Yosemite*



Call 888-646-2244
for Reservations

[Bookmark](#) | [Invite a Friend](#) | [Sign up](#) | [Contact](#) | [Directions](#)

- HOME
 - ACCOMMODATIONS
 - AMENITIES
 - TRAVEL GROUPS
 - SPECIALS & PACKAGES
 - ABOUT YOSEMITE
- RESERVATIONS**
- Arrival:

Must-See Attractions

More Information: [About Yosemite](#) [Attractions](#) [Activities](#) [Entertainment](#) [Shopping](#) [Dining](#)

Exciting Attractions near Yosemite Miner's Inn Hotel

Birdwatching

Yosemite is home to variety of birds, including:

- | | | |
|--------------------|-----------------------|---------------------|
| Stellar's jay | Raven | Great gray owl |
| American robin | Black-headed grosbeak | Peregrine falcon |
| Brewer's blackbird | Red-wing blackbird | Pileated woodpecker |
| Acorn woodpecker | American dipper | Northern goshawk |

Query: *What are must-see attractions at Beijing*



酒店预订 目的地指南 风景图库 列车时刻表 旅游论坛HOT

Hotel ads

预订北京酒店
一方订房网
订房专线 400-819-1189

五星酒店 四星酒店 三星酒店 二星酒店

- 北京亚洲大酒店 ★★★★★ ¥1050
- 北京京都信苑饭店 ★★★★★ ¥750
- 强强(北京)国际商务酒店 ★★★★☆ ¥458
- 北京京仪大酒店 ☆☆☆☆☆ ¥680
- 北京大悦城酒店公寓 ☆☆☆☆☆ ¥788
- 北京融金国际酒店 ☆☆☆☆☆ ¥570
- 北京凯莱大酒店 ★★★★★ ¥550
- 北京宝辰饭店 ☆☆☆☆☆ ¥458
- 北京亮马河大厦 ★★★★★ ¥738
- 北京华威商务全套房酒店 ☆☆☆☆☆ ¥588
- 北京西单美爵酒店 ☆☆☆☆☆ ¥690
- 北京金桥国际公寓 ☆☆☆☆☆ ¥468
- 北京美华世纪国际酒店 ☆☆☆☆☆ ¥588
- 北京清华紫光国际交流中心 ★★★★★ ¥450
- 北京瑞银特公寓酒店 ☆☆☆☆☆ ¥418
- 北京万丰世纪国际大酒店 ☆☆☆☆☆ ¥248

目的地旅游指南 - 直辖市旅游指南 - 北京旅游指南

北京旅游景点 重庆旅游景点 上海旅游景点 天津旅游景点

- 北京旅游指南 - 北京旅游景点 - 北京游记攻略 - 北京特产美食 - 北京当地资讯 - 北京风景美图 - 北京酒店特惠 -

详细的北京景点,北京旅游景点介绍为您到北京旅游提供旅游帮助

推荐阅读

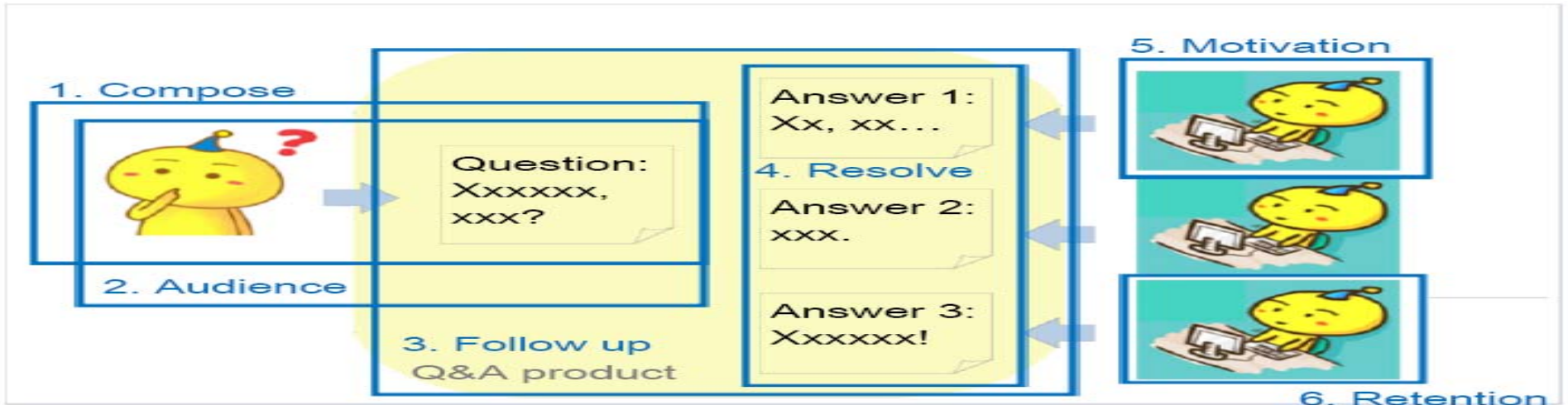
- 北京旅游地图
- 北京首都博物馆
- 制造艳遇 北京美女出没地点大全
- 北京鸟巢
- 北京:五大烤鸭经典餐厅全攻略
- 北京北海公园
- 深秋枫叶渐红 北京赏枫攻略
- 北京水立方
- 北京自助游实用省钱之攻略
- 北京欢乐谷
- 北京毛主席纪念堂

北京旅游景点

人文古迹,自然景观,公园游乐场

- 北京首都博物馆
- 北京欢乐谷
- 北京天安门
- 北京焦庄户地道战遗址纪念馆
- 北京五棵松体育馆
- 北京密云黑龙潭
- 北京圣米厄尔教堂
- 北京水立方
- 北京北海公园
- 中国科学技术馆
- 北京八大处公园
- 北京大学
- 北京烟袋斜街
- 北京鸟巢
- 北京毛主席纪念堂
- 北京陶然亭公园
- 北京中央广播电视塔
- 北京密云水库
- 北京仙栖洞
- 北京自然博物馆

Key ML Subroutines of Confucius



- Trigger a question session during search
- Given a question, provide labels for easy organization
- Given a question, find similar questions and their answers
- Evaluate user credentials in a domain sensitive way
- Given a question, route it to domain experts
- Evaluate quality of answers to a question
- Machine-generated answers

Naive User Evaluation

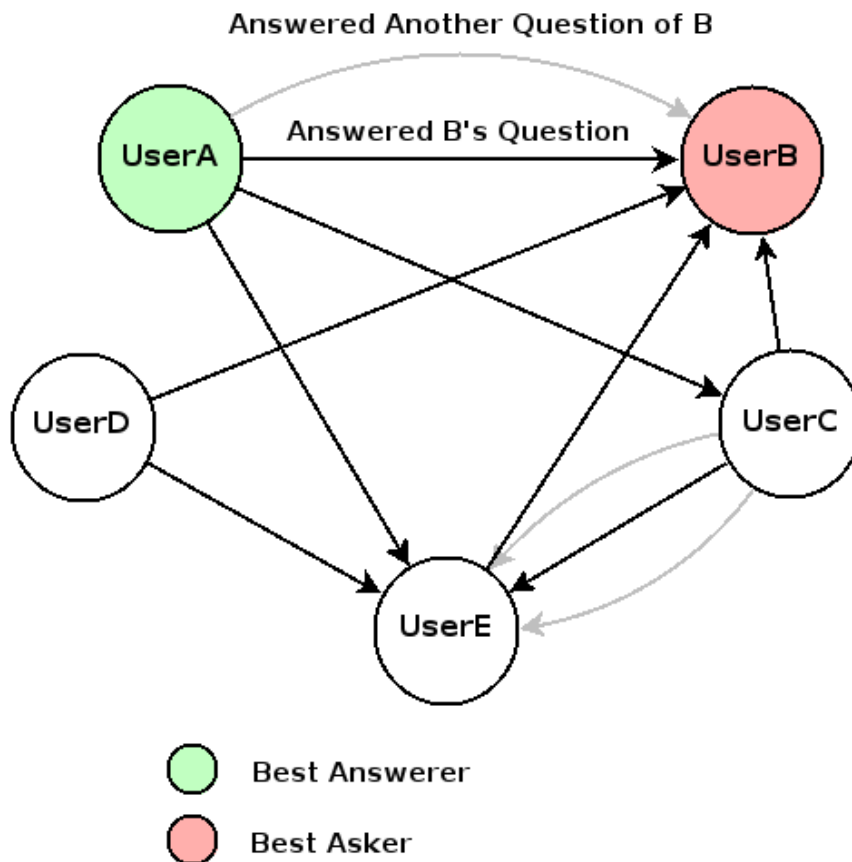
Point system based on *hand-crafted rules*:

- registration → +100 points
- each time login → +5 points
- ask one question → +bonus points
- ask one question → +2 points
- vote on one answer → +1 points
- best answers → +bonus points
- ...

Shortcomings

- **Easily Spammed**
 - Mutual enforcement, answer “friends” questions
 - 1,000 IDs of the same person
 - Copy & paste others’ answers
 - Advertising posts
- **Freshness**
 - User's recent activities are not emphasized

Link-based User Credential Ranking: HITS



QA pairs \rightarrow User Relation
Ranking user using HITS*

| | In links | Out Links |
|---|----------|-----------|
| A | 0 | 4 (3) |
| B | 4 | 0 |
| C | 1 | 4 (2) |
| D | 0 | 2 |
| E | 5 (3) | 1 |

* HITS is based on Zoltan et al, *Questioning Yahoo Answers*. QAWeb, WWW2008

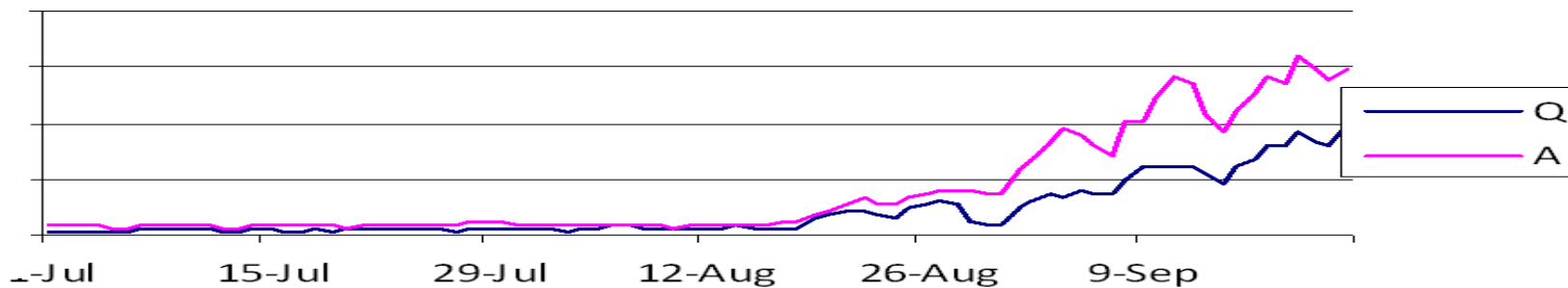
Q&A/Blog/BBS Search

- Lack of links
- Links can be easily spammed
- User credential can help ranking

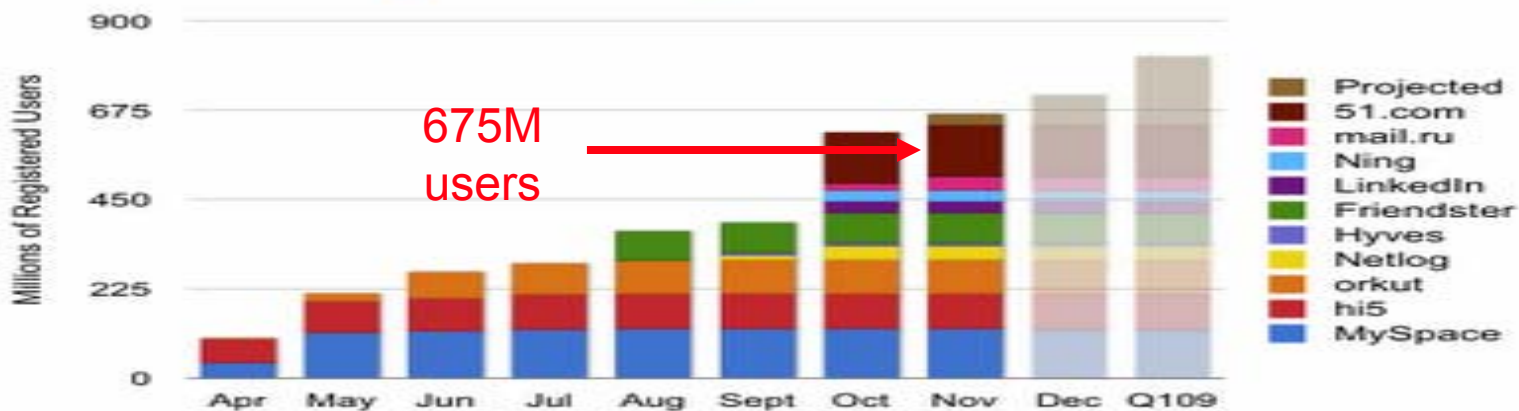
| Q | QA pair ranking |
|----|-----------------|
| A1 | 0.7 |
| A2 | 0.2 |
| A3 | 0.9 |

Data Mining Impact & Opportunities

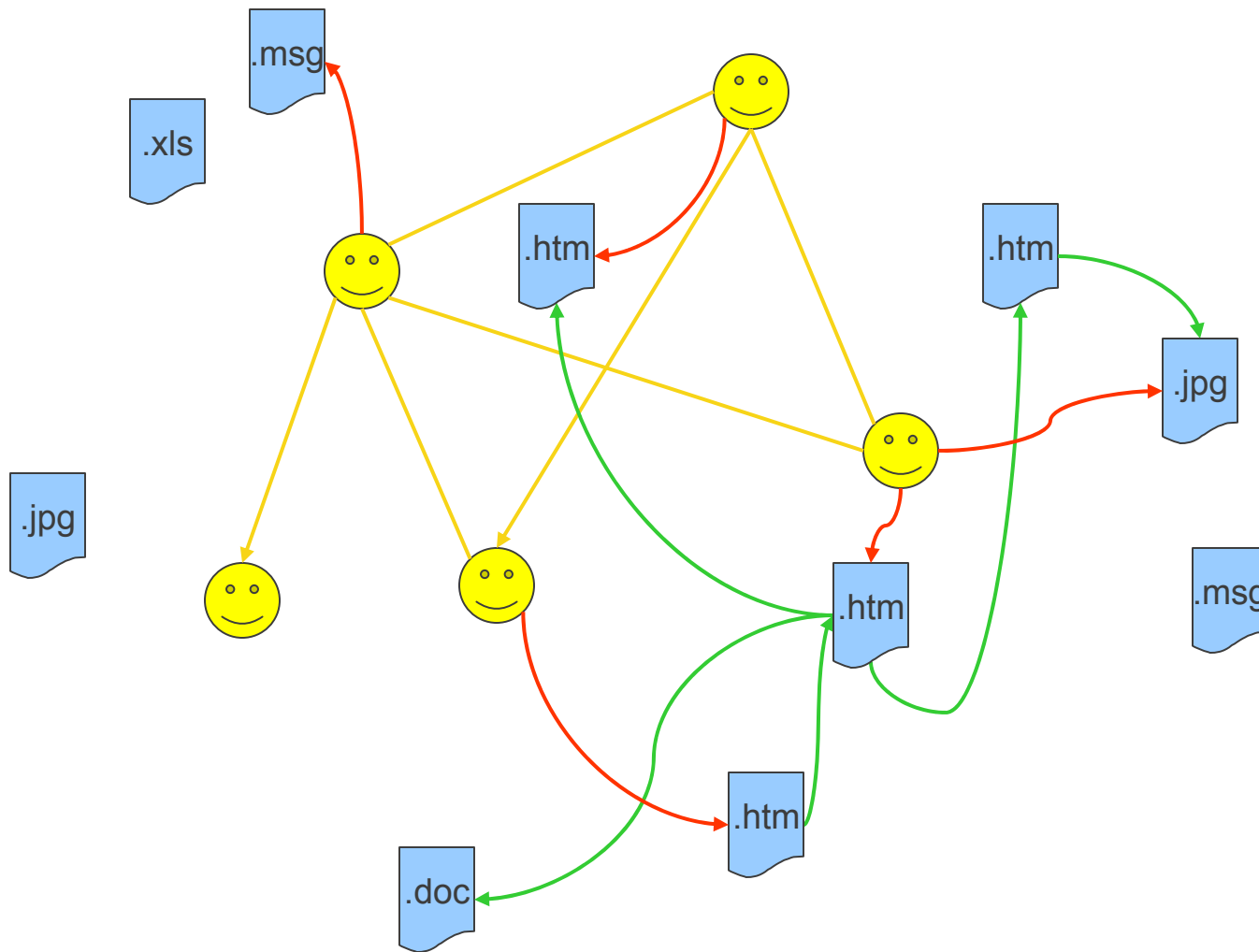
Confucius Growth



opensocial reach

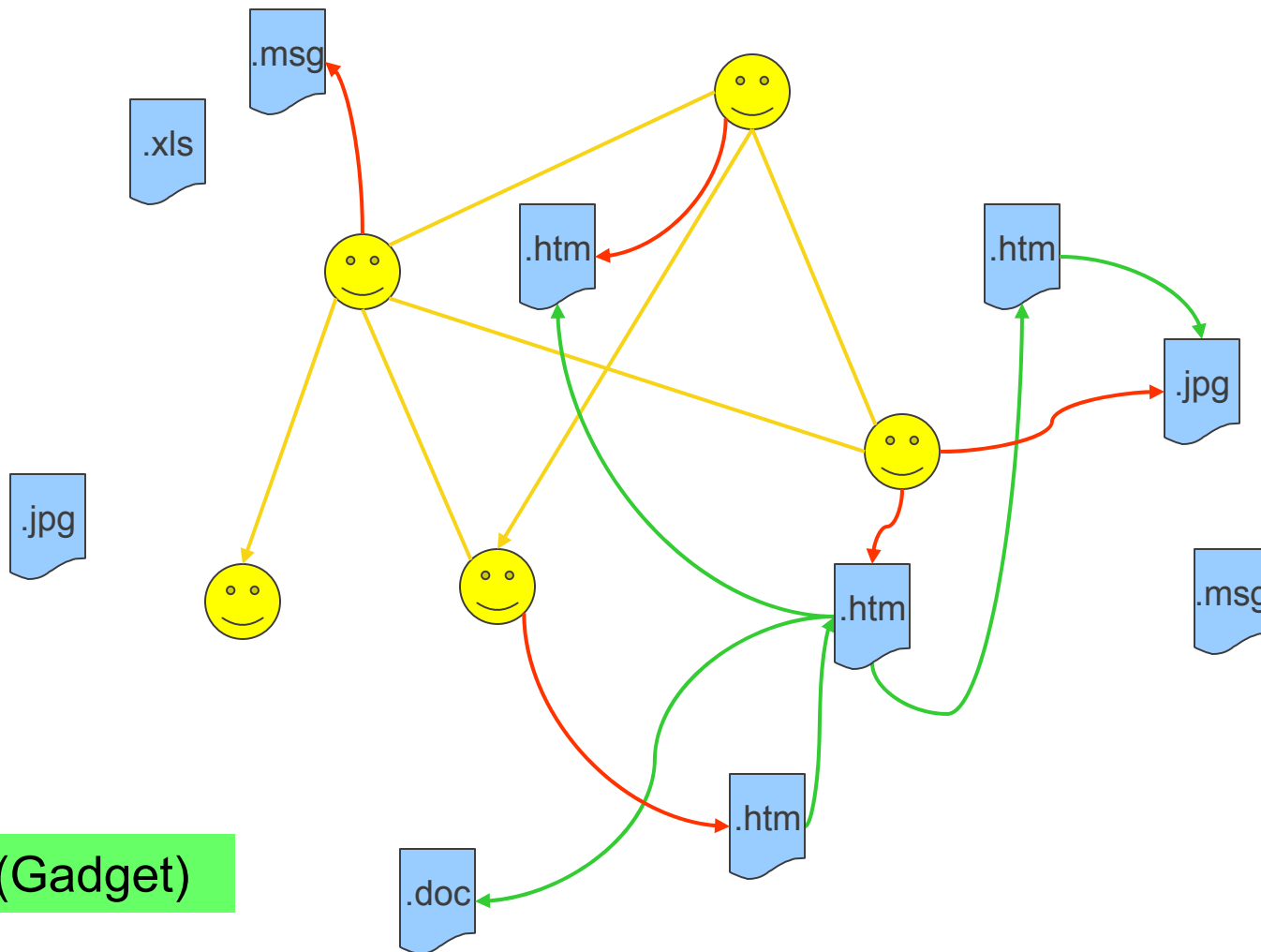


Web 2.0 --- Web with People



+ Social Platforms

App (Gadget)

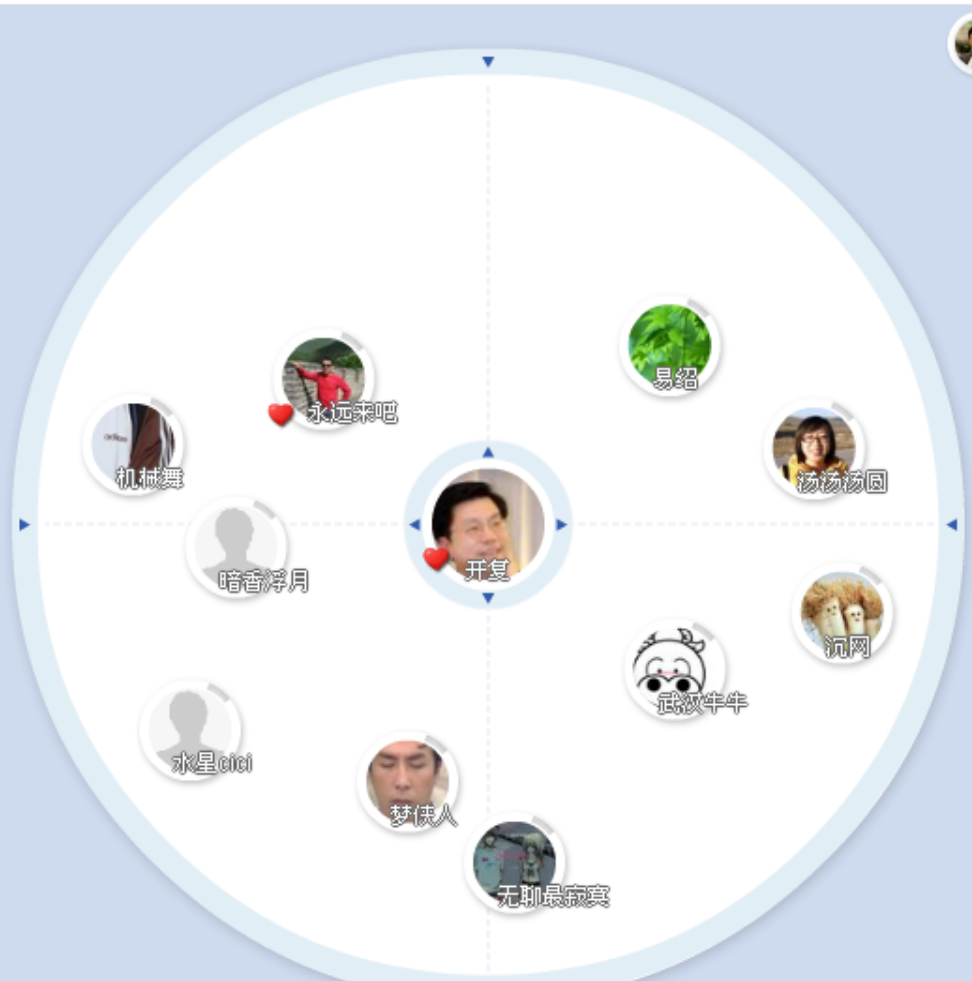


App (Gadget)

我的朋友圈 | 我的朋友

上一步

回到自己



邀请朋友加入来吧

邮件

发送邀请

我要群发邀请 >

看看我的朋友在哪



我的朋友圈 | 我的朋友

◀ 上一步 | ▶

回到自己

永远来吧 (离线) ❤ 我的好友 ✕

批量上传照片!

男 37岁 北京 ➡

📧 📧 📧 📧

[夜幕诱惑,诱的就是你!...](#) [08-1-3]

[这有没GOOGLE公司的伙...](#) [07-8-20]

[白领的家庭聚会](#) [07-12-30]

[白领的家庭聚会](#) [07-12-30]

[白领的家庭聚会](#) [07-12-30]

✕ 移除好友

邀请朋友加入来吧

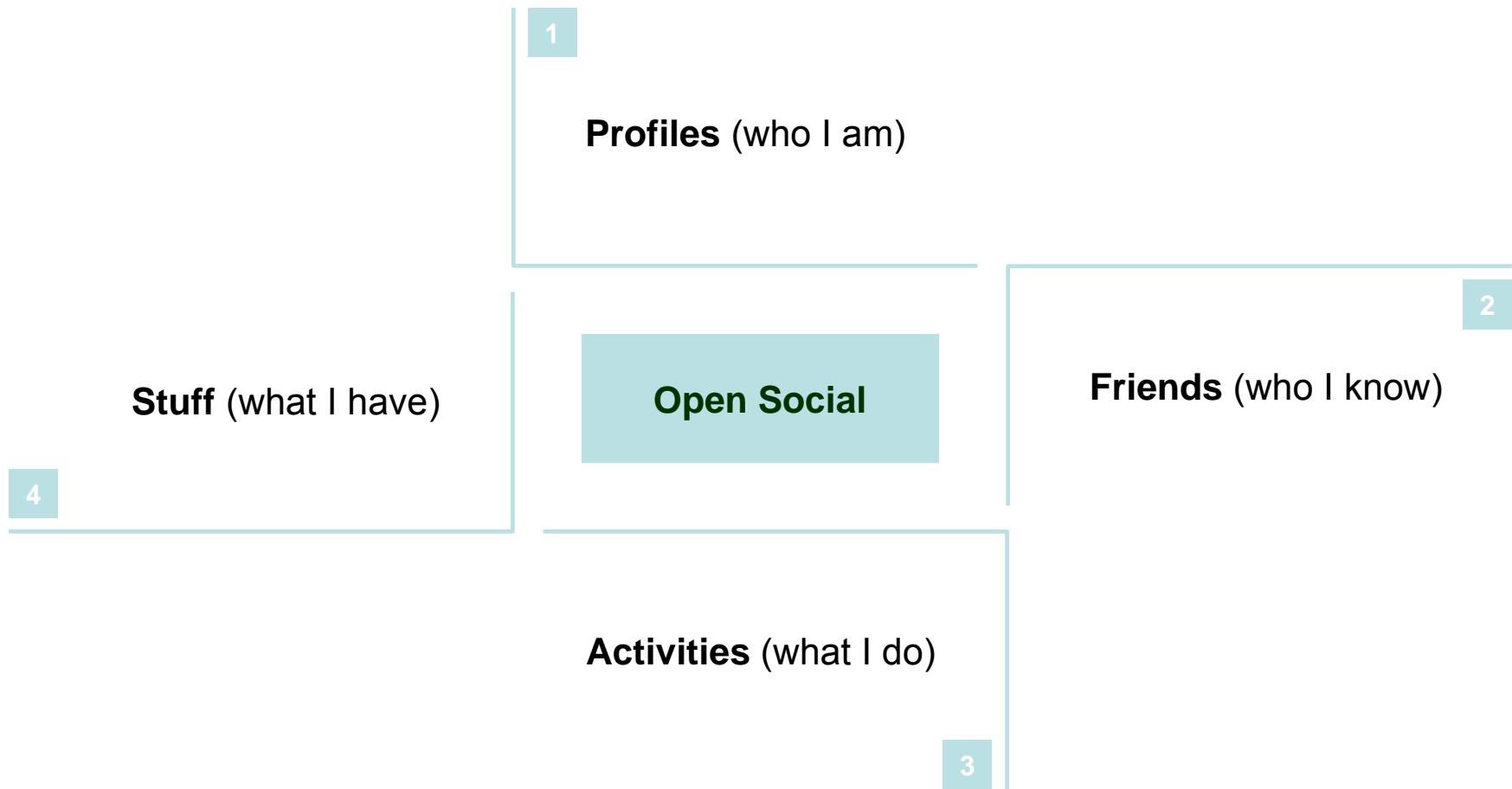
邮件

[我要群发邀请>](#)

看看我的朋友在哪



What users are interested in?



between religiousness of a country and its progress.

I love India. It bothers me, therefore, that the society seems to be much more religious now than in 1960s when I was growing up there. Yes, India has made good progress and possibly economically India is in the best situation now than any other time in the last 200 years. But when one thinks at what has happened in many other Asian countries (not only China, but also in Korea, Taiwan, Singapore, Malaysia, and soon in Vietnam) and compares current Indian situation to what could be, it becomes very depressing. And sitting in this wonderful lounge at Beijing Airport, and comparing this to the lounges in the Mumbai or Delhi Airports, this thought is obvious.

[TECHNICAL THOUGHTS](#), [GENERAL UPDATES](#) | [1 COMMENT](#) »

Beijing Trip

Posted by Ramesh on December 7th, 2008

The last 3 days I have been in Beijing to attend SKG2008. I was requested to give a keynote talk at this conference — Semantics, Knowledge, and Grid — and I talked about the EventWeb ideas.

Though I came here only after about 7 months, this trip showed me a bit more of how rapidly China is transformed. It does not feel like a developing country — all the facilities and the infrastructure makes it look better than many developed countries. Of course, people tell me that once you go away from a few top places like Beijing and Shanghai, the story is different. Even if that is the case, what China has accomplished seems to be unparalleled in the history. Being Indian, it is natural for me to think about India and I feel very depressed about India. In fact I feel worried a bit even about USA

Archives

- [December 2008](#)
- [November 2008](#)
- [October 2008](#)
- [September 2008](#)
- [August 2008](#)
- [July 2008](#)
- [June 2008](#)
- [May 2008](#)
- [April 2008](#)
- [March 2008](#)
- [February 2008](#)
- [January 2008](#)
- [December 2007](#)
- [November 2007](#)
- [October 2007](#)
- [September 2007](#)
- [August 2007](#)
- [July 2007](#)
- [June 2007](#)
- [May 2007](#)
- [April 2007](#)
- [March 2007](#)
- [February 2007](#)
- [January 2007](#)
- [December 2006](#)
- [November 2006](#)
- [October 2006](#)
- [September 2006](#)
- [August 2006](#)
- [July 2006](#)

来吧主页 > 我的相册

我的主页 资料 朋友 来吧 帖子 相册 日记 礼物 ^{New!} 评价 留言簿

上传照片

您已经使用了1024MB中的2MB (0%)

我的相册(8) 我的相片收藏



北京研究会 (18)
点击 1014 2008-4-29



北京过年 ... (6)
点击 1395 2008-2-13



天涯谷歌会议 (12)
点击 1405 2008-1-28



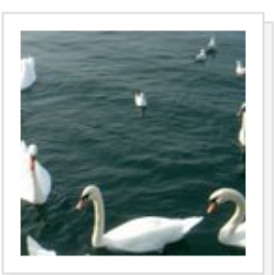
三亚 Goo... (12)
点击 1414 2008-1-19



绿色网络生活 (4)
点击 1331 2007-12-29



成都 De... (6)
点击 1361 2007-12-29



Europe ... (4)
点击 2864 2007-10-24

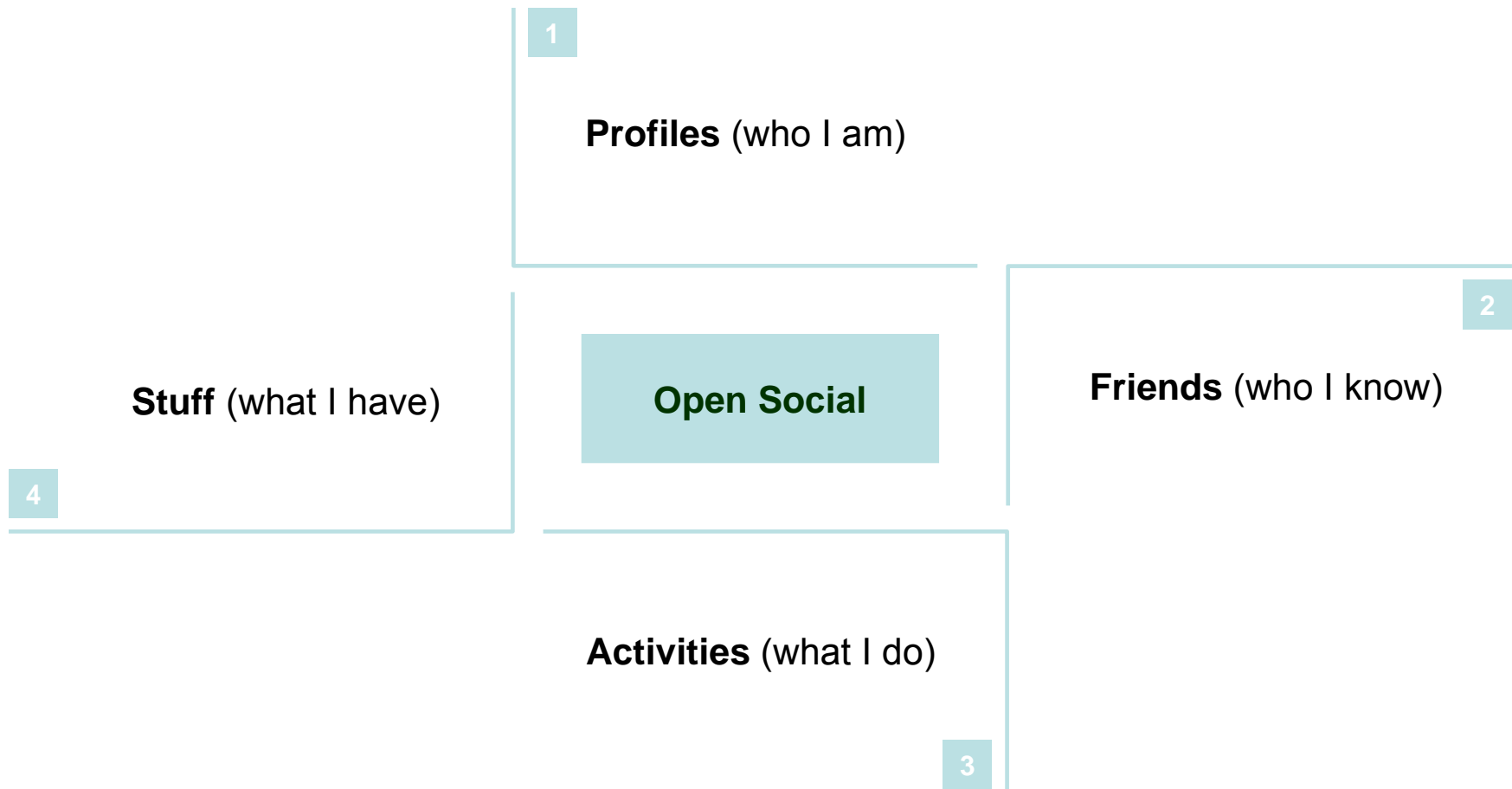


默认相册 (0)
点击 1343

我朋友的相册

-  panjie...的相册
18照片
-  calla_qu的相册
1296照片
-  磨刀美眉的相册
2照片
-  天涯游的相册
85照片
-  三皮儿的相册
61照片
-  钟巍巍的相册
682照片
-  01liux...的相册
75照片
-  efatao...的相册
9照片
-  开复的相册

Open Social APIs



邮件 *

姓名

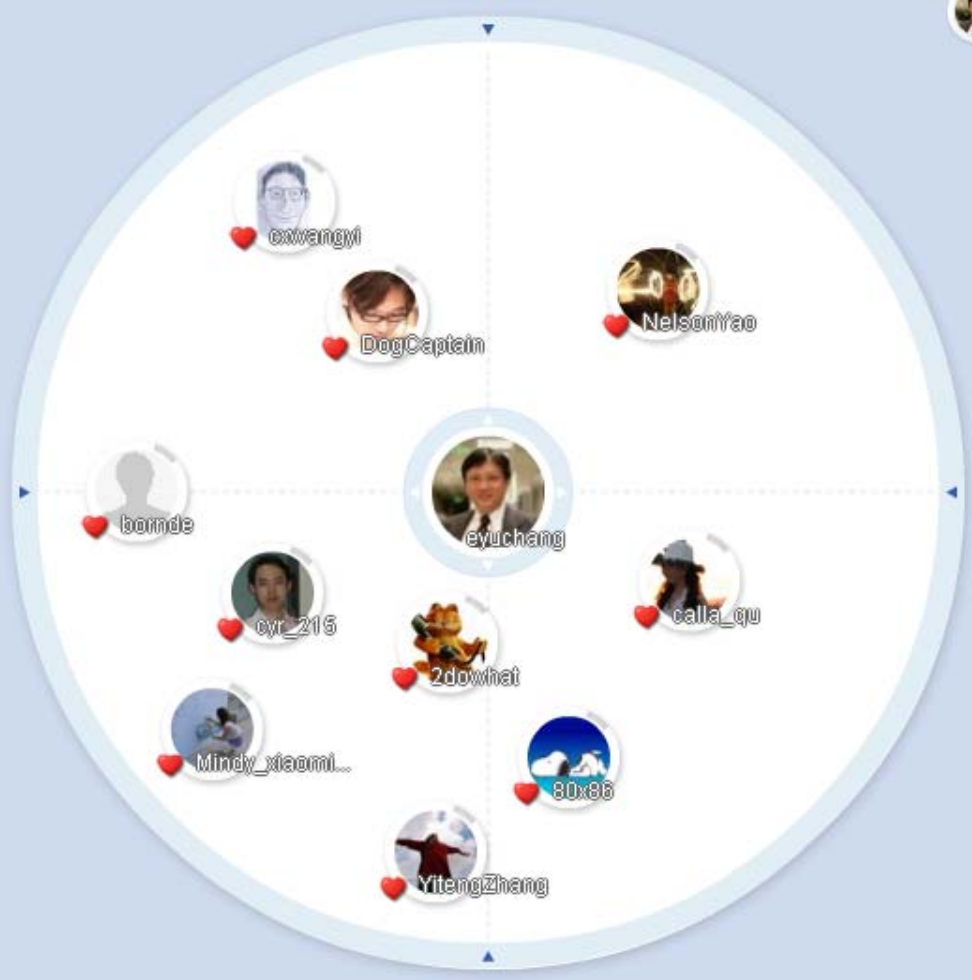
[我要群发邀请 >](#)



我的朋友圈 | 我的朋友

上一步

回到自己



Google™





OpenSocial



Personalized Search Example

- Infer relevance through social networks
- Query “fuji” can return
 - Fuji mountain
 - Fuji apples
 - Fuji cameras



fuji

Search Images

Search the Web

[Advanced Image Search](#)
[Preferences](#)

Moderate SafeSearch is on

Images Showing: All image sizes

Try your search on [Yahoo](#), [Ask](#), [AllTheWeb](#), [Live](#), [PicSearch](#), [Ditto](#), [Getty](#), [Creatas](#), [FreeFoto](#), [WebShots](#), [NASA](#), [Flickr](#), [deviantART](#), [Photobucke](#)



Mt Fuji, Japan
1572 x 1069 - 414k - jpg



Mount Fuji
800 x 639 - 100k - jpg



Northwestern view of Mt. Fuji over ...
And here is the Mount Fuji that the ..
...





fuji

Search Images

Search the Web

[Advanced Image Search](#)
[Preferences](#)

[Moderate SafeSearch is on](#)

Images Showing:

Try your search on [Yahoo](#), [Ask](#), [AllTheWeb](#), [Live](#), [PicSearch](#), [Ditto](#), [Getty](#), [Creatas](#), [FreeFoto](#), [WebShots](#), [NASA](#), [Flickr](#), [deviantART](#), [Photobuck](#)



(Apples, Fuji) Fuji apples are an

...

765 x 792 - 37k - jpg

www.all-creatures.org



fuji apple

300 x 294 - 17k - jpg

www.wisegeek.com



Organic - **Apples, Fuji**

375 x 375 - 67k - jpg

www.cleanfoodconnection.com



fuji apple Manufacturer

800 x 600 - 81k - jpg

www.supplierlist.com



fuji

Search Images

Search the Web

[Advanced Image Search](#)
[Preferences](#)

Moderate SafeSearch is on

Images Showing: All image sizes

Try your search on [Yahoo](#), [Ask](#), [AllTheWeb](#), [Live](#), [PicSearch](#), [Ditto](#), [Getty](#), [Creatas](#), [FreeFoto](#), [WebShots](#), [NASA](#), [Flickr](#), [deviantART](#), [Photobucket](#)



... "as is typical of Fuji cameras ... **fujifilm digital camera**, digital, ...
400 x 400 - 78k - jpg
www.livingroom.org.au



fujifilm digital camera, digital, ...
464 x 254 - 13k - jpg
www.fujifilm-cameras.com



Fuji cameras, one with face ...
425 x 313 - 35k - jpg
www.gadgetell.com



Fuji fujifilm finepix A800
425 x 290 - 34k - jpg
www.gadgetell.com

我的朋友圈 | 我的朋友

上一步

回到自己



ifilm digital camera, digital, ...
464 x 254 - 13k - jpg
www.fujifilm-cameras.com



Fuji cameras, one with face ...
425 x 313 - 35k - jpg
www.gadgetell.com



Fuji Camera
450 x 333 - 440k - bmp
emergencygadget.com



cheap fuji digital camera
400 x 318 - 14k - jpg
cheap-digital-camera.com.au

Google™

邀请朋友加入来吧

邮件 *
姓名

发送邀请

我要群发邀请 >

看看我的朋友在哪



Personalized Recommendation



Recommendation Systems

- Photo/Video Recommendation
- Friend Recommendation
- Community/Forum Recommendation
- Ads Matching

- Performance Requirements
 - Scalability, scalability, scalability

Outline

- Applications
 - Confucius
 - OpenSocial
- Key Subroutines for Mining Massive SNS
 - Clustering [\[ECML 08\]](#)
 - Frequent Itemset Mining [\[ACM RS 08\]](#)
 - Combinational Collaborative Filtering [\[KDD 08\]](#)
 - with PLSA
 - with LDA
 - Support Vector Machines [\[NIPS 07\]](#)
- Distributed Computing Perspectives

Outline

- Applications
 - Confucius
 - OpenSocial
- Key Subroutines
 - Clustering [\[ECML 08\]](#)
 - Frequent Itemset Mining (FIM)
 - Combinational Collaborative Filtering
 - with PLSA
 - with LDA
 - Support Vector Machines
- Distributed Computing Perspectives

Task: Targeting Ads at SNS Users

Users

| | | | | | |
|---|--|--|--|---|--|
|  | <p>miss_ming 女 发消息 282 关注 0</p> |  | <p>宝贝玛德莲 女 发消息 12 关注 3 回复 0</p> |  | <p>蛋笑笑 女 发消息 7424 关注 33</p> |
|  | <p>combaby秋千闲逛 发消息 1494 关注 2</p> |  | <p>桃花临水 3.16 命中注定我 发消息 575 关注 51</p> |  | <p>trovbley GOLD VS LEAF 发消息 81 关注 16 回复 1</p> |
|  | <p>WTN 男 发消息 540 关注 0</p> |  | <p>ys5354 男 发消息 268 关注 2</p> |  | <p>诺百成 男 发消息 571 关注 4</p> |
|  | <p>32679319 男 发消息 569 关注 259</p> |  | <p>famously 女, 22岁 发消息 567 关注 73</p> |  | <p>飞天鼠标 男, 19岁, 河南 发消息 979 关注 112</p> |

Ads



December 12th, 08

nd Seal

Mining Profiles, Friends & Activities for Relevance

我的资料

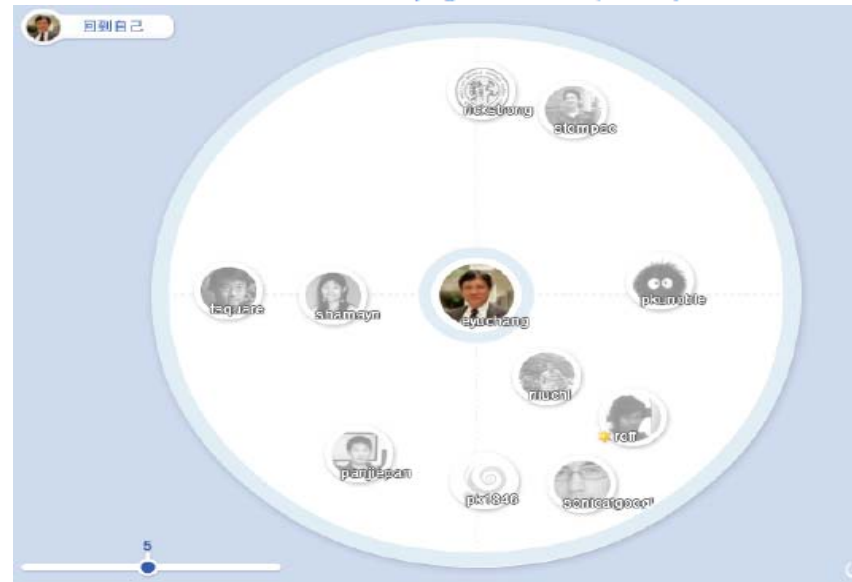


姓名: eyuchang
 真实姓名: 张智威
 性别: 男
 星座: 狮子
 住址: 北京
 家乡: 甘肃
 大学: Stanford
 公司: Google
 书籍: The Castle (Franz Kafka)
 The Brothers Karamazov (Fyodor Dostoevsky)
 Essays of Friedrich Schiller
 Iphigenia in Tauris (Goethe)

登录: 2008年9月23日
 人气: 7556次访问
 积分: 8777
 好友: 88
 照片: 62
 帖子: 10

相册

| | | | | | |
|---|-----------------------------|---|--------------------------------------|---|----------------------------------|
|  | 北京研究会 2008-4-29 10照片 |  | 北京过年 2008 2008-2-13 6照片 |  | 大连谷歌会议 2008-1-28 12照片 |
|  | 绿色网络生活 2007-12-29 4照片 |  | 成都 December ... 2007-12-29 6照片 |  | Europe Trip 2007-10-24 4照片 |



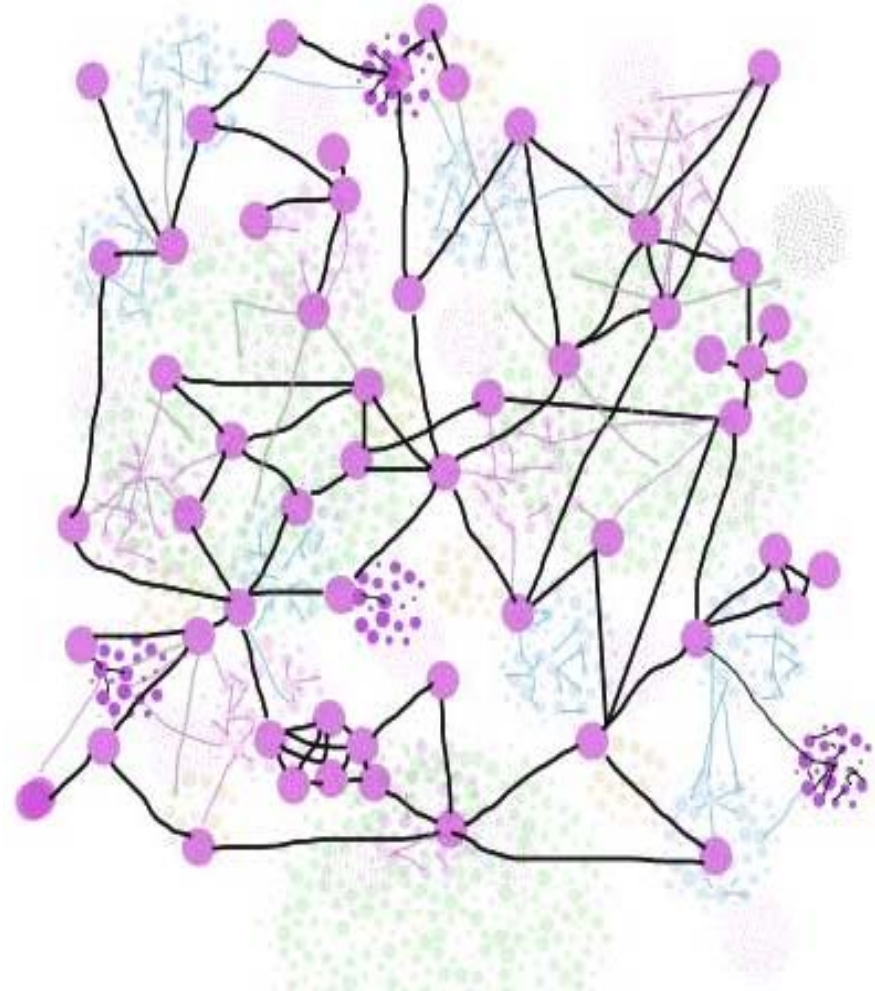
帖子

标题

- [余建漫画] 小狐狸KIKO的QQ表情下载
- [杨欣] 波霸杨欣激情
- [体操] 莫慧兰筹备退役选手就业辅导基金 关注无名选手
- [张梓琳] 中国张梓琳获世界小姐冠军全过程回放
- [摄影爱好者] 兵马俑在大英博物馆
- [浪漫韩剧] 最新搜集文根英图集
- [谣言谎报] 外电称西门子中国有近一半的业务涉及行贿

Consider also User *Influence*

- Advertisers consider users who are
 - Relevant
 - *Influential*
- SNS Influence Analysis
 - Centrality
 - Credential
 - Activeness
 - etc.



Outline

- Emerging Applications
 - Social networks
 - Personalized Information retrieval
- Key Subroutines
 - Clustering [\[ECML 08\]](#)
 - Frequent Itemset Mining (FIM)
 - Combinational Collaborative Filtering
 - with PLSA
 - with LDA
 - Support Vector Machines

Collaborative Filtering

Based on *membership* so far,
and *memberships* of others



Predict further *membership*

Photos/Videos

| | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | 1 | 1 | 1 | | | | | | |
| | 1 | | 1 | 1 | | 1 | | 1 | | 1 |
| | | | | | 1 | | 1 | | | 1 |
| | 1 | | 1 | | 1 | 1 | | | | |
| | | 1 | | | | | | | | |
| | | | | | | 1 | 1 | | | |
| | | | 1 | | | | | 1 | | |
| 1 | 1 | | | | | | | | | |
| | 1 | | | | | | | | 1 | |
| 1 | | | | | | | | | | 1 |
| | 1 | 1 | 1 | 1 | 1 | | | | | |

Users

Some Queries

Based on *partially*
observed matrix



Predict *unobserved* entries



I. Will user *i* enjoy photo *j*?

II. Will user *i* be interesting to user *j*?

III. Will photo *i* be related to photo *j*?

Photos/Videos

| | | | | | | | | | | | |
|--|---|---|---|---|---|---|---|---|---|---|--|
| | ? | | 1 | 1 | 1 | | ? | | | | |
| | 1 | ? | 1 | 1 | | 1 | | 1 | | 1 | |
| | | | | | 1 | ? | 1 | | | 1 | |
| | 1 | | 1 | ? | 1 | 1 | | | | | |
| | | 1 | | | | | ? | | | | |
| | ? | | | | | 1 | 1 | | | | |
| | | | 1 | | | | | 1 | ? | | |
| | 1 | 1 | | | ? | | | | | | |
| | | 1 | | | | ? | | | 1 | | |
| | 1 | | | | | | | ? | | 1 | |
| | | 1 | 1 | 1 | 1 | 1 | ? | | | | |

Users

FIM-based Recommendation



To grow the base, we need association rules

- An association rule: $a, b, c \longrightarrow d$
- A Bayesian interpretation: $P(d | a, b, c) = \frac{N(a, b, c, d)}{N(a, b, c)}$
- The key is to count the occurrences (*support*) of itemsets $N(\dots)$

FIM Preliminaries

- Observation 1: If an item A is not frequent, any pattern contains A won't be frequent [R. Agrawal]
→ use a threshold to eliminate infrequent items
 ~~$\{A\}$~~ → ~~$\{A, B\}$~~
- Observation 2: Patterns containing A are subsets of (or found from) transactions containing A [J. Han]
→ divide-and-conquer: select transactions containing A to form a conditional database (CDB), and find patterns containing A from that conditional database
 $\{A, B\}, \{A, C\}, \{A\} \rightarrow \text{CDB } A$
 $\{A, B\}, \{B, C\} \rightarrow \text{CDB } B$
- Observation 3: Some patterns may be found in multiple CDBs

Preprocessing

f a c d g i m p

a b c f l m o

b f h j o

b c k s p

a f c e l p m n

f: 4

c: 4

a: 3

b: 3

m: 3

p: 3

o: 2

d: 1

e: 1

g: 1

h: 1

i: 1

k: 1

l: 1

n: 1

f c a m p

f c a b m

f b

c b p

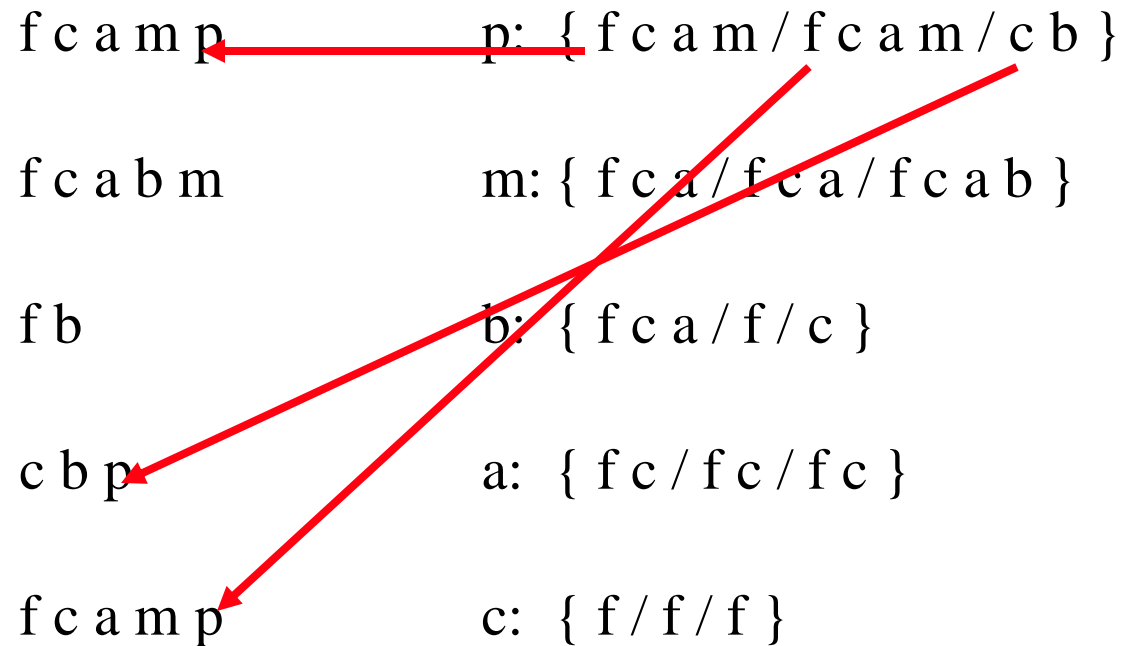
f c a m p

- According to Observation 1, we count the support of each item by scanning the database, and eliminate those infrequent items from the transactions.
- Sort items in each transaction by the order of descending support value.

Parallel Projection

- According to Observation 2, we construct CDB of item A ; then from this CDB, we find those patterns containing A
- How to construct the CDB of A ?
 - If a transaction contains A , this transaction should appear in the CDB of A
 - Given a transaction $\{B, A, C\}$, it should appear in the CDB of A , the CDB of B , and the CDB of C
- Dedup solution: using the order of items:
 - sort $\{B, A, C\}$ by the order of items $\rightarrow \langle A, B, C \rangle$
 - Put $\langle \rangle$ into the CDB of A
 - Put $\langle A \rangle$ into the CDB of B
 - Put $\langle A, B \rangle$ into the CDB of C

Example of Projection

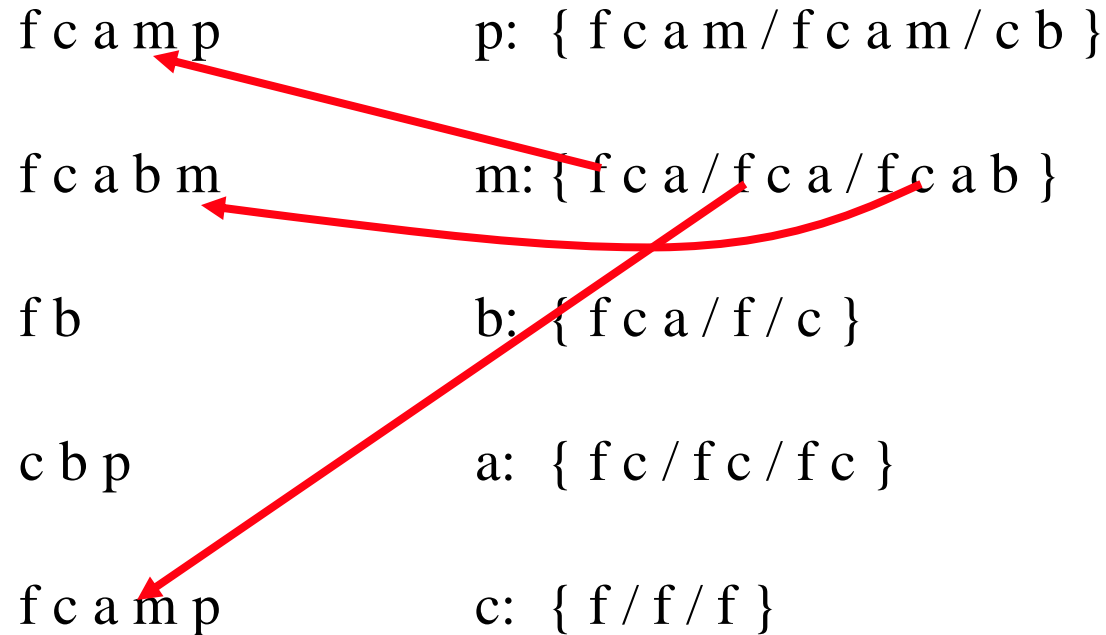


Example of Projection of a database into CDBs.

Left: sorted transactions in order of *f*, *c*, *a*, *b*, *m*, *p*

Right: conditional databases of frequent items

Example of Projection



Example of Projection of a database into CDBs.

Left: sorted transactions;

Right: conditional databases of frequent items

Example of Projection

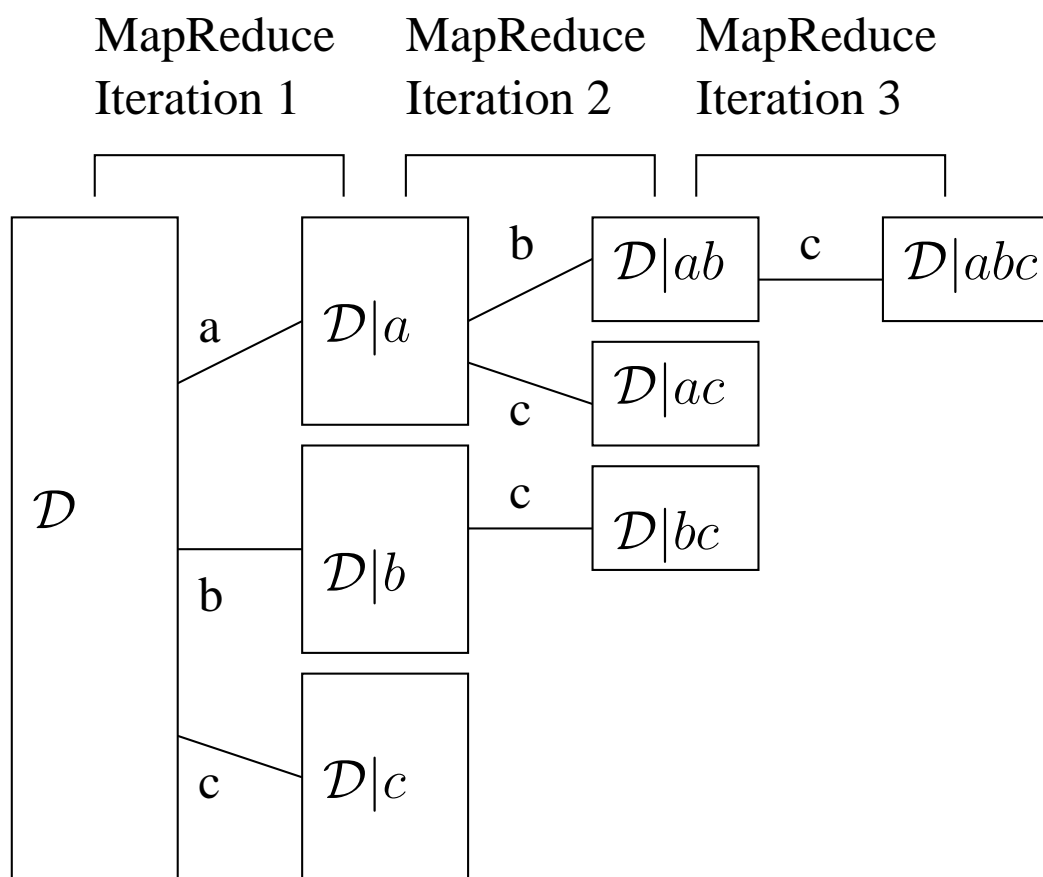
| | |
|-----------|--------------------------------|
| f c a m p | p: { f c a m / f c a m / c b } |
| f c a b m | m: { f c a / f c a / f c a b } |
| f b | b: { f c a / f / c } |
| c b p | a: { f c / f c / f c } |
| f c a m p | c: { f / f / f } |

Example of Projection of a database into CDBs.

Left: sorted transactions;

Right: conditional databases of frequent items

Recursive Projections [H. Li, et al. ACM RS]



- Recursive projection form a search tree
- Each node is a CDB
- Using the order of items to prevent duplicated CDBs.
- Each level of breath-first search of the tree can be done by a MapReduce iteration.
- Once a CDB is small enough to fit in memory, we mine this CDB, and no more growth of the sub-tree.

Projection using MapReduce

| Map inputs (transactions) key="": value | Sorted transactions (with infrequent items eliminated) | Map outputs (conditional transactions) key: value | Reduce inputs (conditional databases) key: value | Reduce outputs (patterns and supports) key: value | |
|---|--|---|--|---|--|
| f a c d g i m p | f c a m p | p: f c a m m: f c a a: f c c: f | p: {fcam/fcam/cb} p:3, pc:3 | | |
| a b c f l m o | f c a b m | m: f c a b b: f c a a: f c c: f | | m f : 3 m c : 3 m a : 3 m f c : 3 m f a : 3 m c a : 3 m f c a : 3 | |
| b f h j o | f b | b: f | | | |
| b c k s p | c b p | p: c b | | b: { f c a / f / c } | b : 3 |
| a f c e l p m n | f c a m p | b: c p: f c a m m: f c a a: f c c: f | | a: { f c / f c / f c } | a : 3 a f : 3 a c : 3 a f c : 3 |
| | | | c: { f / f / f } | c : 3 c f : 3 | |

Outline

- Applications
 - Confucius
 - OpenSocial
- Key Subroutines
 - Clustering
 - Frequent Itemset Mining (FIM)
 - Combinational Collaborative Filtering
 - with PLSA
 - with LDA
 - Support Vector Machines

Collaborative Filtering

Based on *membership* so far,
and *memberships* of others



Predict further *membership*

Forums/Communities

| | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | 1 | 1 | 1 | | | | | | |
| | 1 | | 1 | 1 | | 1 | | 1 | | 1 |
| | | | | | 1 | | 1 | | | 1 |
| | 1 | | 1 | | 1 | 1 | | | | |
| | | 1 | | | | | | | | |
| | | | | | | 1 | 1 | | | |
| | | | 1 | | | | | 1 | | |
| 1 | 1 | | | | | | | | | |
| | 1 | | | | | | | | 1 | |
| 1 | | | | | | | | | | 1 |
| | 1 | 1 | 1 | 1 | 1 | | | | | |

Users

Collaborative Filtering

Based on *membership* so far,
and *memberships* of others



Predict further *membership*

Forums/Communities

| | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | 1 | 3 | 1 | | | | | | |
| | 2 | | 1 | 2 | | 1 | | 3 | | 1 |
| | | | | | 1 | | 5 | | | 1 |
| | 5 | | 3 | | 1 | 1 | | | | |
| | | 1 | | | | | | | | |
| | | | | | | 1 | 4 | | | |
| | | | 2 | | | | | 1 | | |
| 1 | 2 | | | | | | | | | |
| | 1 | | | | | | | | 5 | |
| 1 | | | | | | | | | | 1 |
| | 1 | 4 | 1 | 3 | 6 | | | | | |

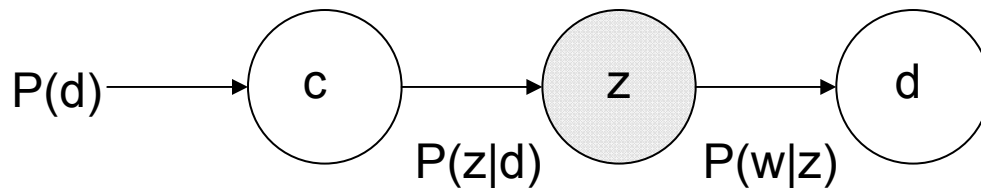
Users

Notations

- Given a collection of co-occurrence data
 - Community: $C = \{c_1, c_2, \dots, c_N\}$
 - User: $U = \{u_1, u_2, \dots, u_M\}$
 - Description: $D = \{d_1, d_2, \dots, d_V\}$
 - Latent aspect: $Z = \{z_1, z_2, \dots, z_K\}$
- Models
 - Baseline models
 - Community-User (C-U) model
 - Community-Description (C-D) model
 - CCF: Combinational Collaborative Filtering
 - *Combines* both baseline models

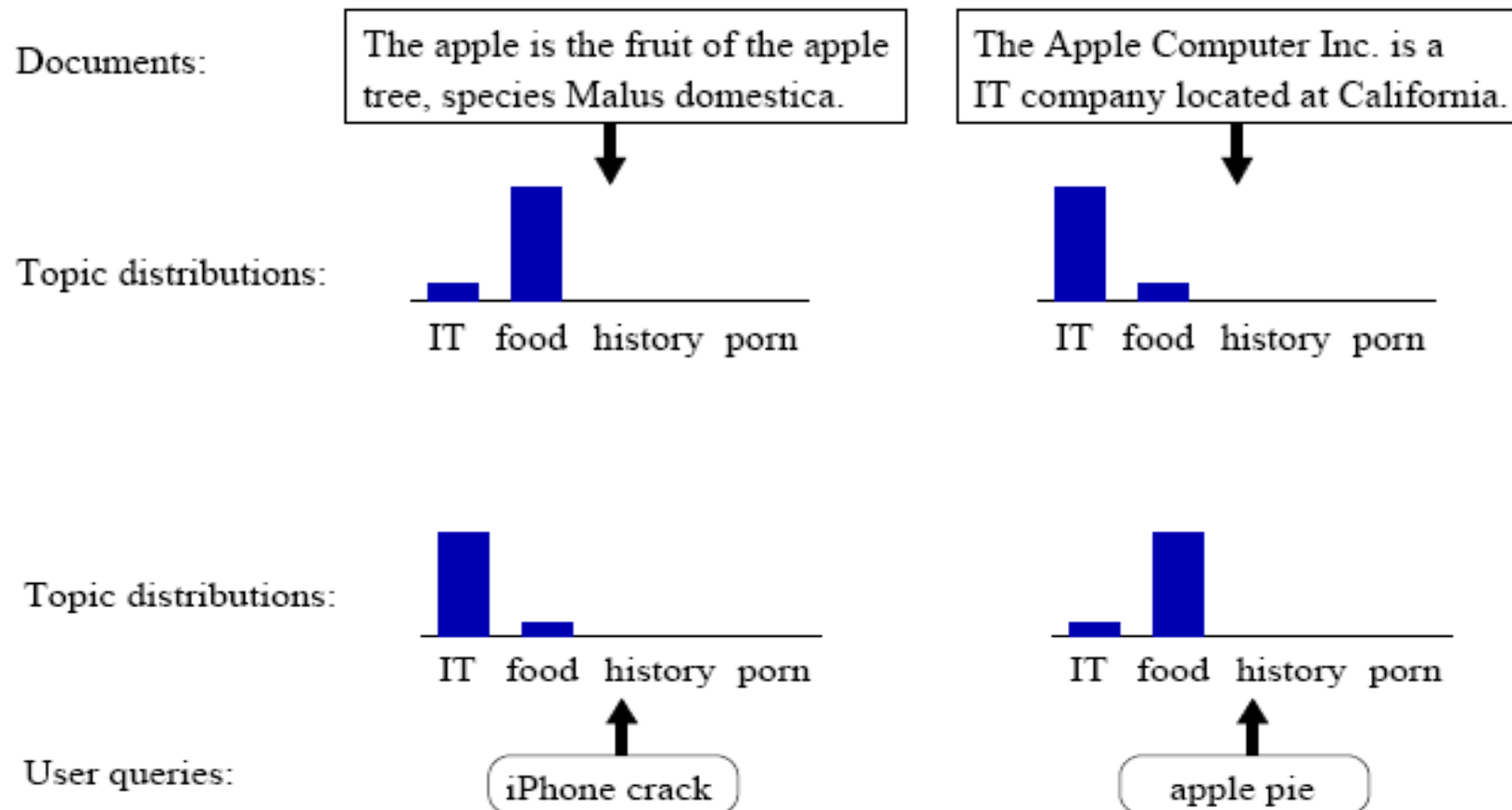
Probabilistic Latent Semantic Analysis (PLSA) [Hoffman 1999; Hoffman 2004]

- Document is viewed as a bag of words
- A *latent semantic layer* is constructed in between documents and words
- $P(d, c) = P(d|c) P(c) = P(c) \sum_z P(d|z) P(z|c)$



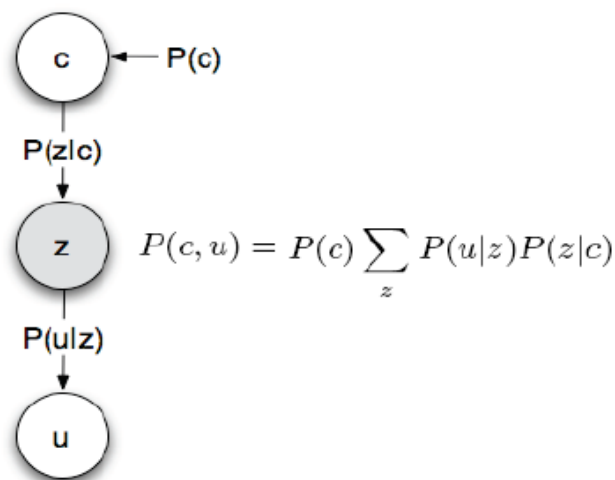
- Probability delivers explicit meaning
 - $P(c|c)$, $P(d|d)$, $P(d, c)$
- Model learning via EM or Gibbs sampling

Example of Latent Analysis



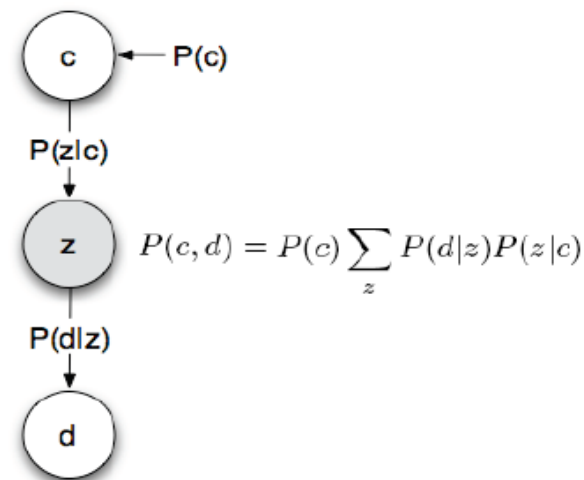
Baseline Models

Community-User (C-U) model



- Community is viewed as a **bag of users**
- **c** and **u** are rendered conditionally independent by introducing **z**
- Generative process, for each user u
 1. A community c is chosen uniformly
 2. A topic z is selected from $P(z|c)$
 3. A user u is generated from $P(u|z)$

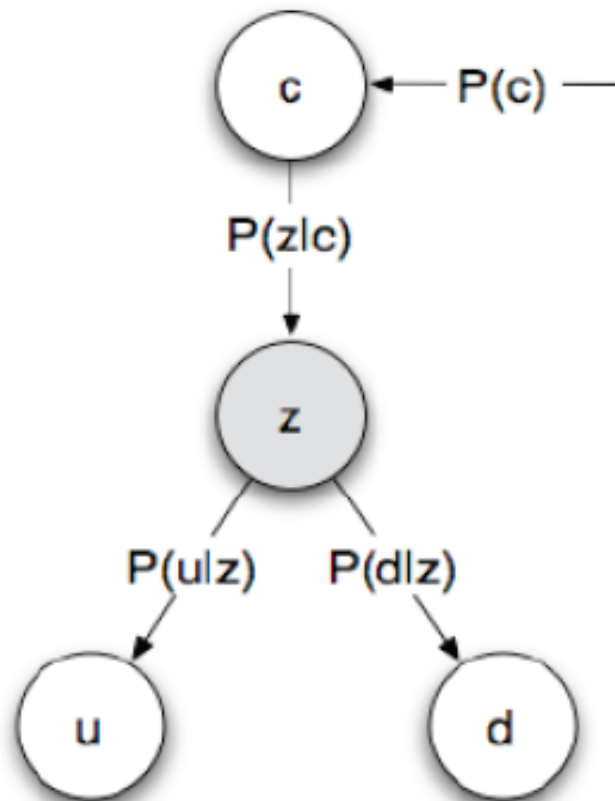
Community-Description (C-D) model



- Community is viewed as a **bag of words**
- **c** and **d** are rendered conditionally independent by introducing **z**
- Generative process, for each word d
 1. A community c is chosen uniformly
 2. A topic z is selected from $P(z|c)$
 3. A word d is generated from $P(d|z)$

CCF Model [Chen, et. al. KDD 08]

Combinational Collaborative Filtering (CCF) model



- CCF *combines* both baseline models
- A community is viewed as
 - a bag of users AND a bag of words
- By adding C-U, CCF can perform personalized recommendation which C-D alone *cannot*
- By adding C-D, CCF can perform better recommendation than C-U alone, which may suffer from sparsity
- CCF can do that C-U and C-D cannot
 - $P(d|u)$, relate user to word
 - Useful for user targeting ads

Empirical Study

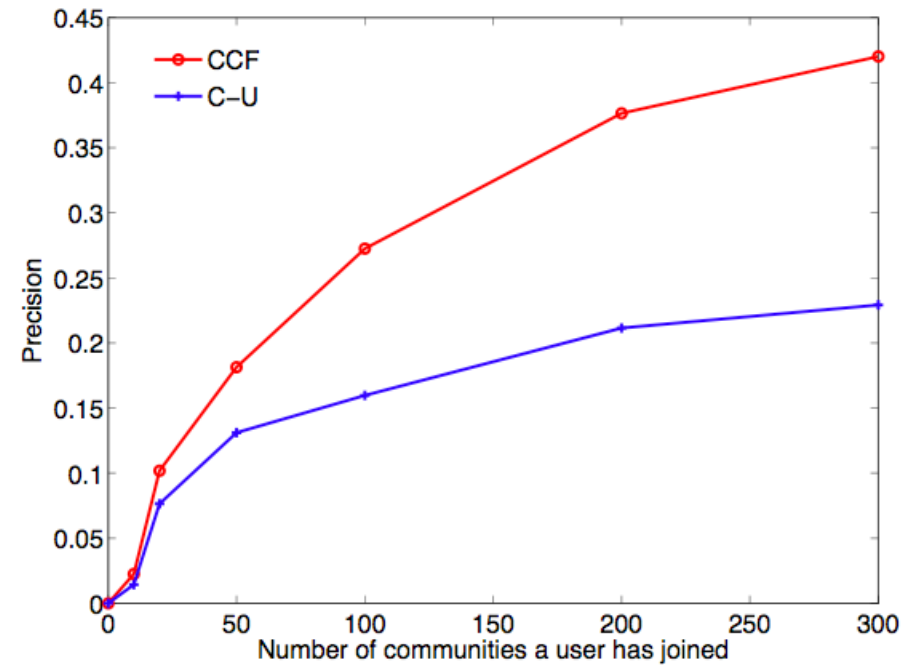
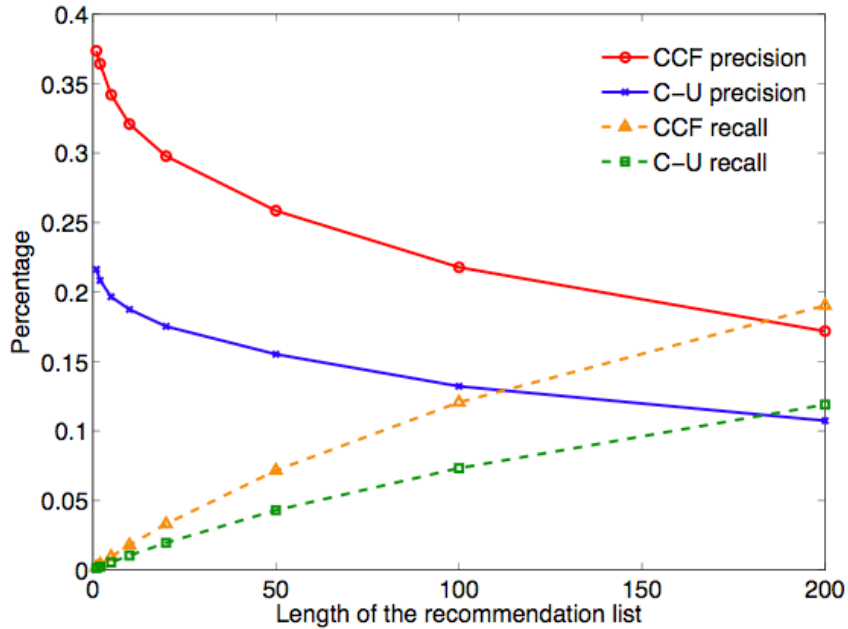
- Orkut Dataset
 - Collected in July, 2007
 - Two types of data were extracted
 - Community-user, community-description
 - 312,385 users
 - 109,987 communities
- Machine farm
 - Up to 200 machines in Google datacenters
 - Each machine is configured with:
 - A CPU faster than 2GHz
 - Memory larger than 4GBytes
- Evaluations
 - Community recommendation
 - Speedup

Community Recommendation

- Evaluation Method
 - Leave-one-out: randomly delete one community for each user
 - Check if a removed community can be recovered?
- Evaluation metric
 - Precision and Recall



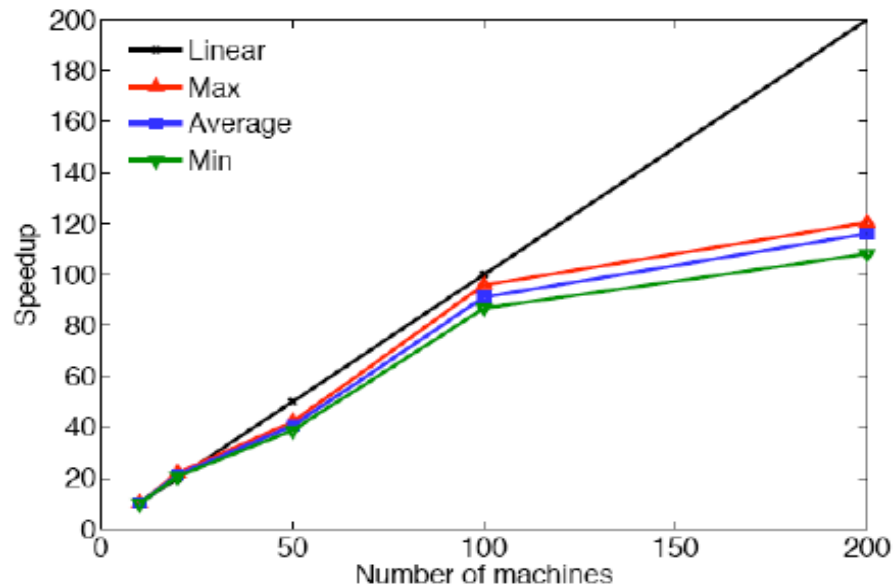
Results



□ CCF outperforms C-U

□ The more information, the higher accuracy

Gibbs Sampling MapReduce Speedup

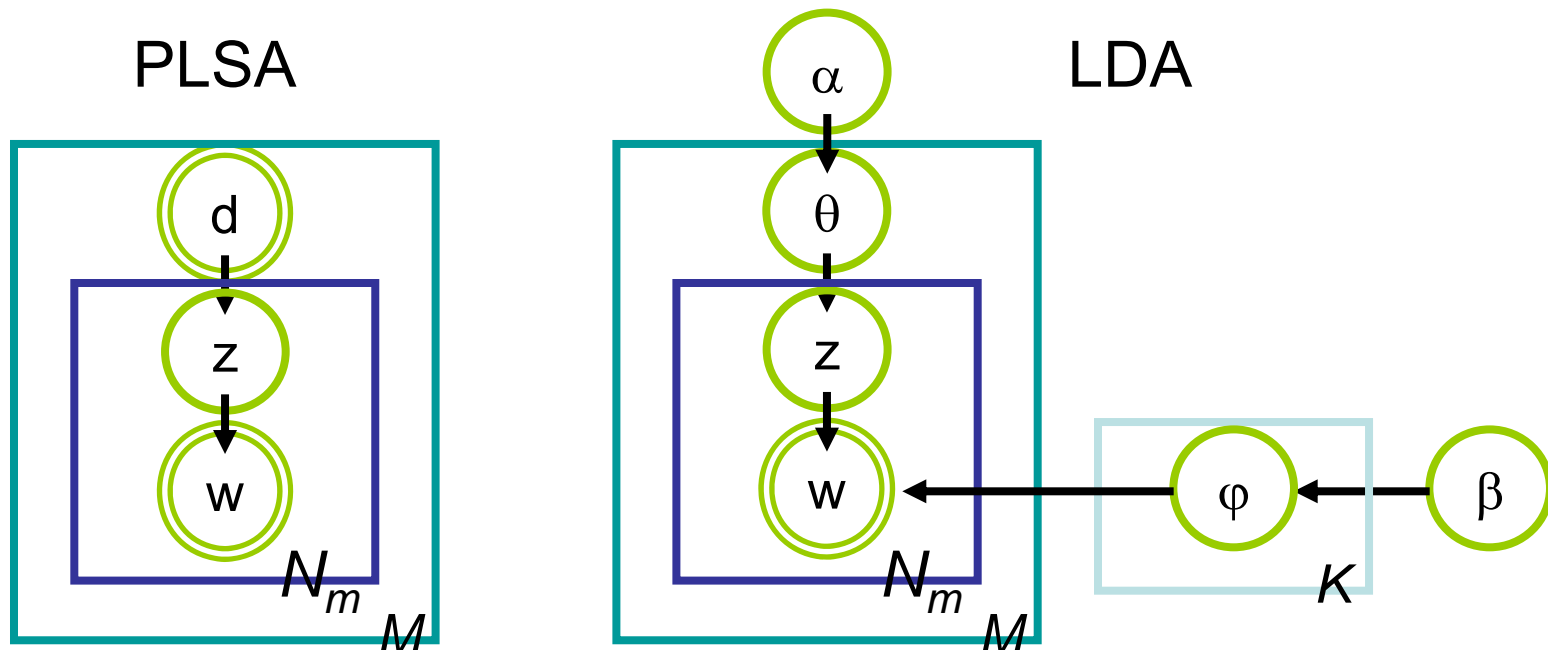


| Machines | Time (sec.) | Speedup |
|----------|-------------|---------|
| 10 | 9,233 | 10 |
| 20 | 4,326 | 21.3 |
| 50 | 2,280 | 40.5 |
| 100 | 1,014 | 91.1 |
| 200 | 796 | 116 |

- The Orkut dataset enjoys a linear speedup when the number of machines is up to 100
- Reduces the training time from one day to less than 14 minutes

Extensions

- Expand CCF to incorporate more types of information
- Replace PLSA with LDA



...Extensions

- Consider time dimension
- Perform incremental learning
- Construct topic hierarchy
- etc...

Outline

- Applications
 - Confucius
 - OpenSocial
- Key Subroutines
 - Clustering
 - Frequent Itemset Mining (FIM)
 - Combinational Collaborative Filtering
 - • with PLSA
 - with LDA

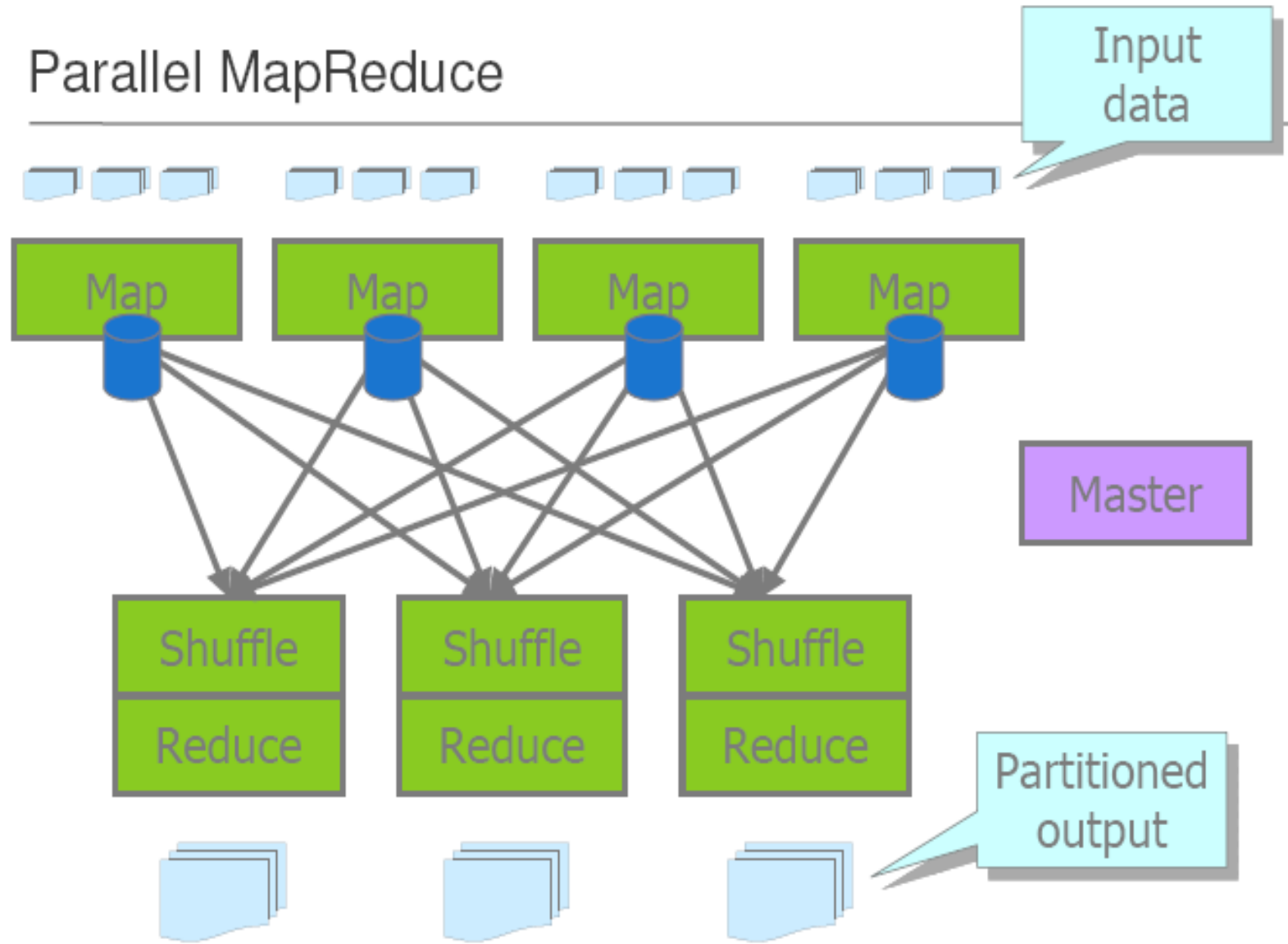
December 12th, 08 NIPS Beyond Search **Support Vector Machines [NIPS 07]**

• Distributed Computing Perspectives

Distributed Computing Perspectives

- **Iterative**
 - Most algorithms do a series of iterations
 - Data dependency: Iteration $t+1$ depends on t
- **Parallelize each iteration**
 - In computation
 - In storage
- **Auto Fault Recovery**
 - Critical for large-scale tasks

Parallel MapReduce

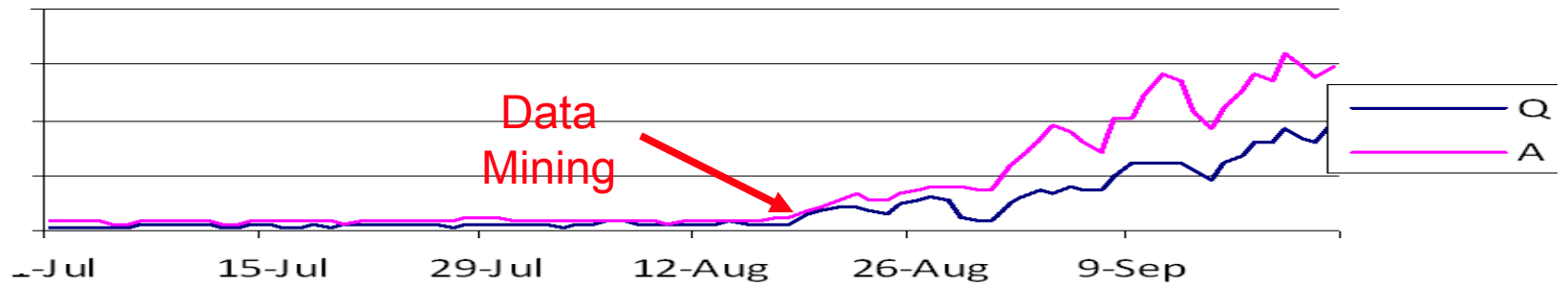


Comparison between Parallel Computing Frameworks

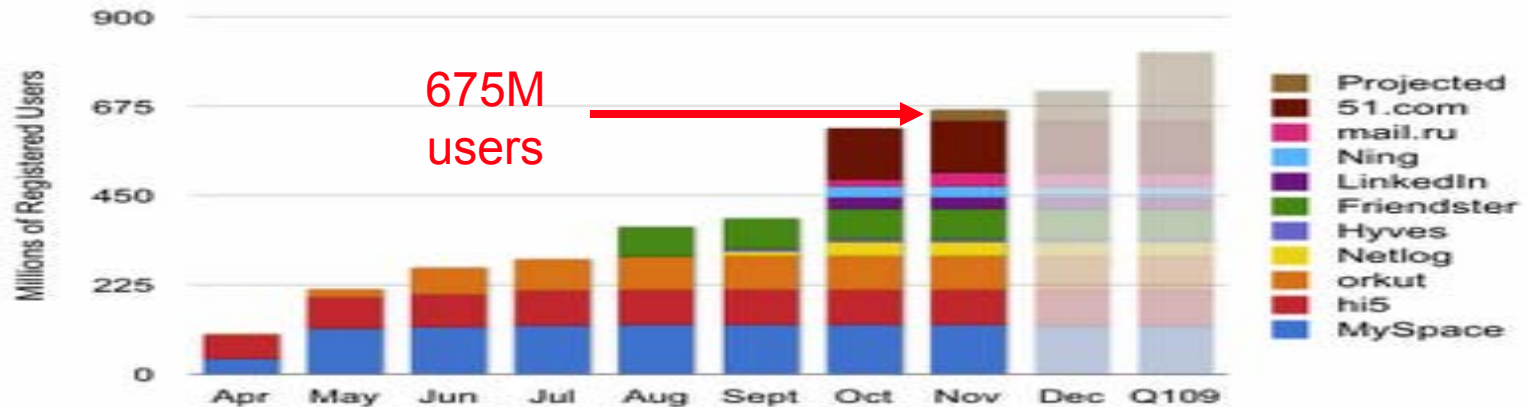
| | MapReduce | Project Doe | MPI |
|--|------------------------|-------------|----------------|
| GFS/IO and task rescheduling overhead between iterations | Yes | No +1 | No +1 |
| Flexibility of computation model | AllReduce only +0.5 | +0.7 | Flexible +1 |
| Efficient AllReduce | Yes +1 | Yes +1 | Yes +1 |
| Recover from faults between iterations | Yes +1 | Yes +1 | No |
| Recover from faults within each iteration | Yes +1 | Yes +1 | No |
| Final Score for scalable machine learning | 3.5 | 4.7 | 3 |

Conclusions...

Confucius Growth



opensocial reach



...Conclusions

- Seven ML subroutines (disciples) of Confucius
- Recommendation is the push model of search
- Recommendation systems demand scalability
- ML algorithms demand “better” distributed computing models than MapReduce

...Conclusions

- Have parallelized key subroutines for mining massive data sets
 - Spectral Clustering [ECML 08]
 - Frequent Itemset Mining [ACM RS 08]
 - Combinational Collaborative Filtering [KDD 08]
 - with PLSA
 - with LDA
 - Support Vector Machines [NIPS 07]
- Relevant papers
 - <http://infolab.stanford.edu/~echang/>
- Open Source PSVM
 - <http://code.google.com/p/psvm/>

References

- [1] Alexa internet. <http://www.alexa.com/>.
- [2] D. M. Blei and M. I. Jordan. Variational methods for the dirichlet process. In Proc. of the 21st international conference on Machine learning, pages 373-380, 2004.
- [3] D. M. Blei, A. Y. Ng, and M. I. Jordan. Latent dirichlet allocation. *Journal of Machine Learning Research*, 3:993-1022, 2003.
- [4] D. Cohn and H. Chang. Learning to probabilistically identify authoritative documents. In Proc. of the Seventeenth International Conference on Machine Learning, pages 167-174, 2000.
- [5] D. Cohn and T. Hofmann. The missing link - a probabilistic model of document content and hypertext connectivity. In *Advances in Neural Information Processing Systems 13*, pages 430-436, 2001.
- [6] S. C. Deerwester, S. T. Dumais, T. K. Landauer, G. W. Furnas, and R. A. Harshman. Indexing by latent semantic analysis. *Journal of the American Society of Information Science*, 41(6):391-407, 1990.
- [7] A. P. Dempster, N. M. Laird, and D. B. Rubin. Maximum likelihood from incomplete data via the em algorithm. *Journal of the Royal Statistical Society. Series B (Methodological)*, 39(1):1-38, 1977.
- [8] S. Geman and D. Geman. Stochastic relaxation, gibbs distributions, and the bayesian restoration of images. *IEEE Transactions on Pattern recognition and Machine Intelligence*, 6:721-741, 1984.
- [9] T. Hofmann. Probabilistic latent semantic indexing. In Proc. of Uncertainty in Artificial Intelligence, pages 289-296, 1999.
- [10] T. Hofmann. Latent semantic models for collaborative filtering. *ACM Transactions on Information System*, 22(1):89-115, 2004.
- [11] A. McCallum, A. Corrada-Emmanuel, and X. Wang. The author-recipient-topic model for topic and role discovery in social networks: Experiments with enron and academic email. Technical report, Computer Science, University of Massachusetts Amherst, 2004.
- [12] D. Newman, A. Asuncion, P. Smyth, and M. Welling. Distributed inference for latent dirichlet allocation. In *Advances in Neural Information Processing Systems 20*, 2007.
- [13] M. Ramoni, P. Sebastiani, and P. Cohen. Bayesian clustering by dynamics. *Machine Learning*, 47(1):91-121, 2002.

References (cont.)

- [14] R. Salakhutdinov, A. Mnih, and G. Hinton. Restricted boltzmann machines for collaborative filtering. In Proc. Of the 24th international conference on Machine learning, pages 791-798, 2007.
- [15] E. Spertus, M. Sahami, and O. Buyukkokten. Evaluating similarity measures: a large-scale study in the orkut social network. In Proc. of the 11th ACM SIGKDD international conference on Knowledge discovery in data mining, pages 678-684, 2005.
- [16] M. Steyvers, P. Smyth, M. Rosen-Zvi, and T. Griffiths. Probabilistic author-topic models for information discovery. In Proc. of the 10th ACM SIGKDD international conference on Knowledge discovery and data mining, pages 306-315, 2004.
- [17] A. Strehl and J. Ghosh. Cluster ensembles - a knowledge reuse framework for combining multiple partitions. Journal on Machine Learning Research (JMLR), 3:583-617, 2002.
- [18] T. Zhang and V. S. Iyengar. Recommender systems using linear classifiers. Journal of Machine Learning Research, 2:313-334, 2002.
- [19] S. Zhong and J. Ghosh. Generative model-based clustering of documents: a comparative study. Knowledge and Information Systems (KAIS), 8:374-384, 2005.
- [20] L. Admic and E. Adar. How to search a social network. 2004
- [21] T.L. Griffiths and M. Steyvers. Finding scientific topics. Proceedings of the National Academy of Sciences, pages 5228-5235, 2004.
- [22] H. Kautz, B. Selman, and M. Shah. Referral Web: Combining social networks and collaborative filtering. Communications of the ACM, 3:63-65, 1997.
- [23] R. Agrawal, T. Imielnski, A. Swami. Mining association rules between sets of items in large databses. SIGMOD Rec., 22:207-116, 1993.
- [24] J. S. Breese, D. Heckerman, and C. Kadie. Empirical analysis of predictive algorithms for collaborative filtering. In Proceedings of the Fourteenth Conference on Uncertainty in Artificial Intelligence, 1998.
- [25] M. Deshpande and G. Karypis. Item-based top-n recommendation algorithms. ACM Trans. Inf. Syst., 22(1):143-177, 2004.

References (cont.)

- [26] B.M. Sarwar, G. Karypis, J.A. Konstan, and J. Reidl. Item-based collaborative filtering recommendation algorithms. In Proceedings of the 10th International World Wide Web Conference, pages 285-295, 2001.
- [27] M.Deshpande and G. Karypis. Item-based top-n recommendation algorithms. ACM Trans. Inf. Syst., 22(1):143-177, 2004.
- [28] B.M. Sarwar, G. Karypis, J.A. Konstan, and J. Reidl. Item-based collaborative filtering recommendation algorithms. In Proceedings of the 10th International World Wide Web Conference, pages 285-295, 2001.
- [29] M. Brand. Fast online svd revisions for lightweight recommender systems. In Proceedings of the 3rd SIAM International Conference on Data Mining, 2003.
- [30] D. Goldberg, D. Nichols, B. Oki and D. Terry. Using collaborative filtering to weave an information tapestry. Communication of ACM 35, 12:61-70, 1992.
- [31] P. Resnik, N. Iacovou, M. Suchak, P. Bergstrom, and J. Riedl. Grouplens: An open architecture for collaborative filtering of netnews. In Proceedings of the ACM, Conference on Computer Supported Cooperative Work. Pages 175-186, 1994.
- [32] J. Konstan, et al. Grouplens: Applying collaborative filtering to usenet news. Communication of ACM 40, 3:77-87, 1997.
- [33] U. Shardanand and P. Maes. Social information filtering: Algorithms for automating “word of mouth”. In Proceedings of ACM CHI, 1:210-217, 1995.
- [34] G. Kinden, B. Smith and J. York. Amazon.com recommendations: item-to-item collaborative filtering. IEEE Internet Computing, 7:76-80, 2003.
- [35] T. Hofmann. Unsupervised learning by probabilistic latent semantic analysis. Machine Learning Journal 42, 1:177-196, 2001.
- [36] T. Hofmann and J. Puzicha. Latent class models for collaborative filtering. In Proceedings of International Joint Conference in Artificial Intelligence, 1999.
- [37] http://www.cs.carleton.edu/cs_comps/0607/recommend/recommender/collaborativefiltering.html
- [38] E. Y. Chang, et. al., Parallelizing Support Vector Machines on Distributed Machines, NIPS, 2007.
- [39] Wen-Yen Chen, Dong Zhang, and E. Y. Chang, Combinational Collaborative Filtering for personalized community recommendation, ACM KDD 2008.
- [40] Y. Sun, W.-Y. Chen, H. Bai, C.-j. Lin, and E. Y. Chang, Parallel Spectral Clustering, ECML 2008.