

# *Integrating Social with Search*

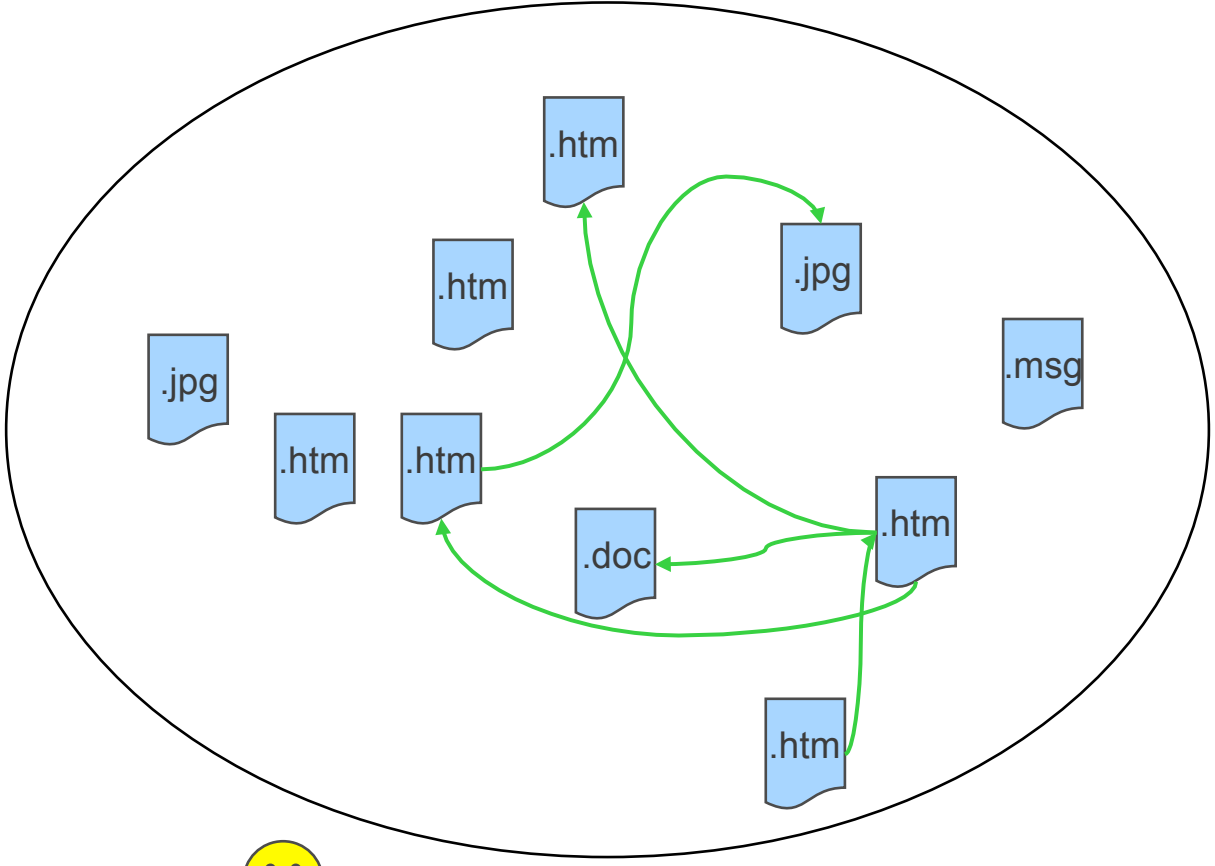
[Edward Chang](#)

*Director, Google Research, Beijing*

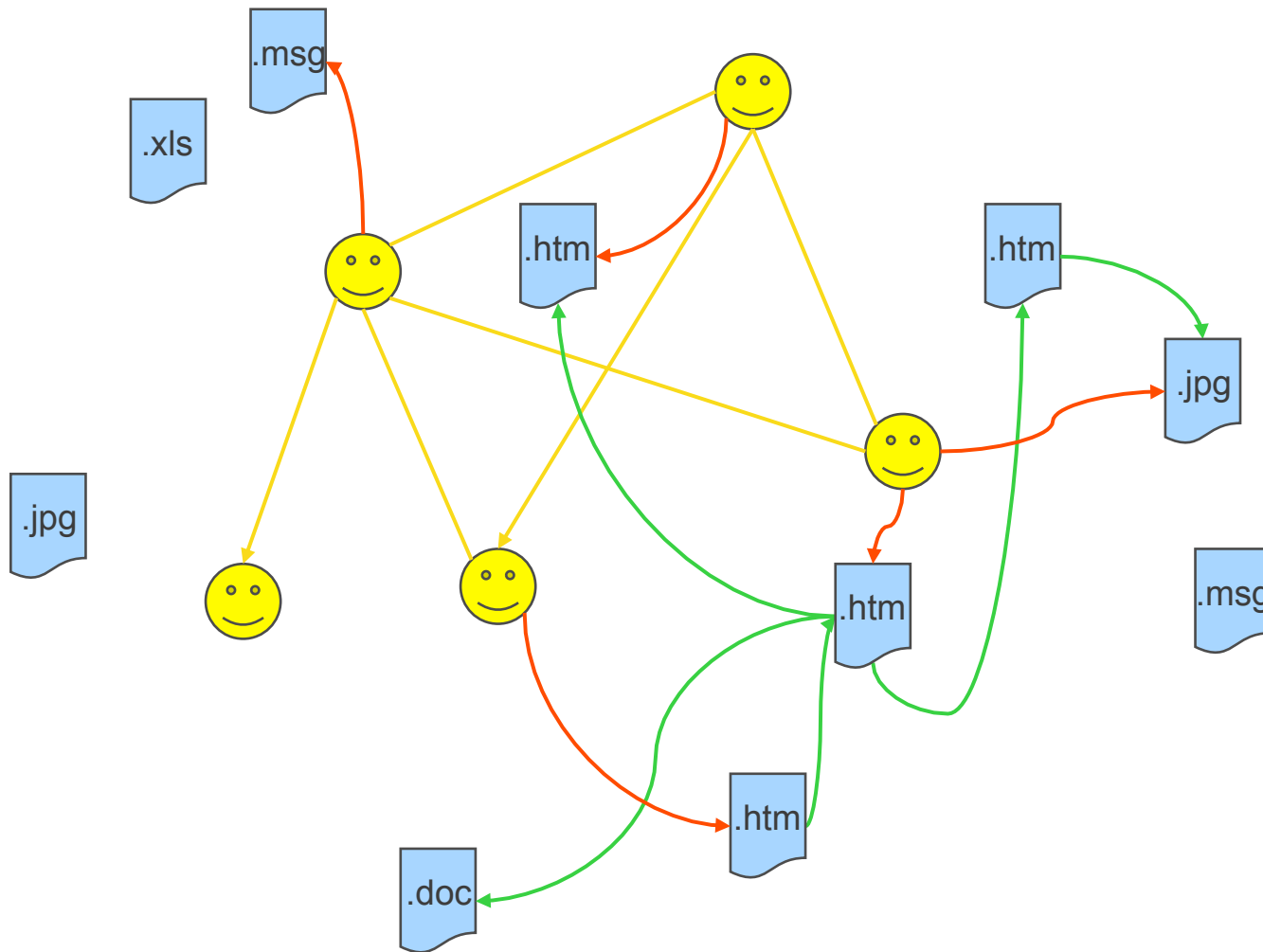
# Related Papers

- **AdHeat (Social Ads):**
  - [AdHeat: An Influence-based Diffusion Model for Propagating Hints to Match Ads](#), H.J. Bao and E. Y. Chang, WWW 2010 (best paper candidate), April 2010.
  - [Parallel Spectral Clustering in Distributed Systems](#), Wen-Yen Chen, Yangqiu Song, Hongjie Bai, Chih-Jen Lin, and E. Y. Chang, IEEE Transactions on Pattern Analysis and Machine Intelligence (PAMI), 2010.
- **UserRank:**
  - Confucius and its Intelligent Disciples, X. Si, E. Y. Chang, Z. Gyongyi, VLDB, September 2010 .
  - Topic-dependent User Rank, Xiance Si, Z. Gyongyi, E. Y. Chang, and M.S. Sun, Google Technical Report.
- **Large-scale Collaborative Filtering:**
  - PLDA+: Parallel Latent Dirichlet Allocation for Large-Scale Applications, ACM Transactions on Internet Technology, 2011.
  - [Collaborative Filtering for Orkut Communities: Discovery of User Latent Behavior](#), W.-Y. Chen, J. Chu, E. Y. Chang, WWW 2009: 681-690.
  - [Combinational Collaborative Filtering for Personalized Community Recommendation](#), W.-Y. Chen, E. Y. Chang, KDD 2008: 115-123.
  - PSVM: Parallelizing SVMs on distributed machines, E. Y. Chang, et al., NIPS 2007.

# Web 1.0



# Web 2.0 --- Web with People



iGoogle

Google Search I'm Feeling Lucky

Advanced Search Language Tools

Show this page every time I start to browse the web. [Make iGoogle my homepage](#)

Change theme from [Spring Scape](#) | [Add stuff](#)

- Home
- YouTube
- Date & Time
- Weather
- Gmail
- CNN.com

- Updates
- Friends

Chat

Search contacts

Ed Chang

革命性创新十年

Weather

Beijing

23°C Current: Light rain  
Wind: E at 26 km/h  
Humidity: 81%

Thu	Fri	Sat	Sun
16°   26°	16°   19°	14°   20°	15°   22°

Date & Time

Thu SEP 16

YouTube

Allow this gadget to:

- Know who I am and see my [Friends](#) group
- Post my activities to Updates.

Save [Learn more](#)

Gmail

Inbox - Compose Mail

Actions...

CNN.com

Pancake house takes on prayer group

Search

Today's Spotlight Videos

Edward Chang Edit My Profile

- News Feed Messages Events Friends (33) Questions More

News Feed Top News · Most Recent 2

Share: Status Question Photo Link Video

Hitesh Gupta http://blogs.timesofindia.indiatimes.com/indus-calling/entry/whose-man-is-that-soldier-fighting-in-kashmir 19 hours ago · Comment · Like Hitesh Gupta likes this. Write a comment...

Ling Ling Mobile Uploads 2 hours ago · Comment · Like · Share 5 people like this. View all 8 comments

Events See All What are you planning?

People You May Know See All Zhao Jing 15 mutual friends Add as friend Dandan Wu 15 mutual friends Add as friend

Requests See All 7 friend requests 26 friend suggestions 1 christmas cheer request

Questions See All Where are some good places to go freshwater fishing in Southern California? Is there a difference between tipping cash and tipping on a credit card?

# Outline


- Search + Social Synergy
- Search → Social
- Social → Search
- Scalability

# Google Q&A (Confucius)

- Developed from 2007 till now @ Beijing
- Launched in more than 60 countries
  - Russia
  - HK
  - Southeast Asia
  - Arab World
  - Sub-Saharan Africa (Baraza)



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

   [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*

THE MINERS INN

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    |

Done

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



## Hotel ads

风景图库 列车时刻表 旅游论坛HOT

**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酒店 ★★★★★ ¥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

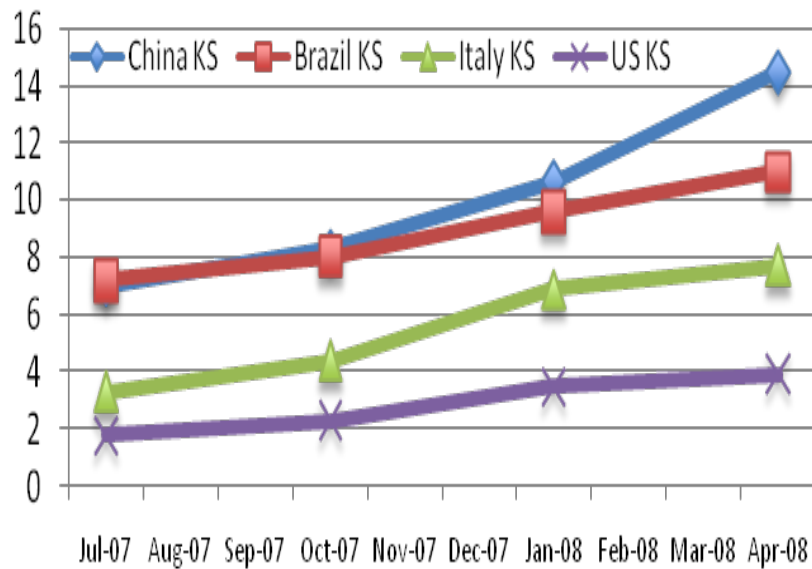
人文古迹,自然景观,公园游乐场

- 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
- 北京园子
- 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
- 北京皇城墙

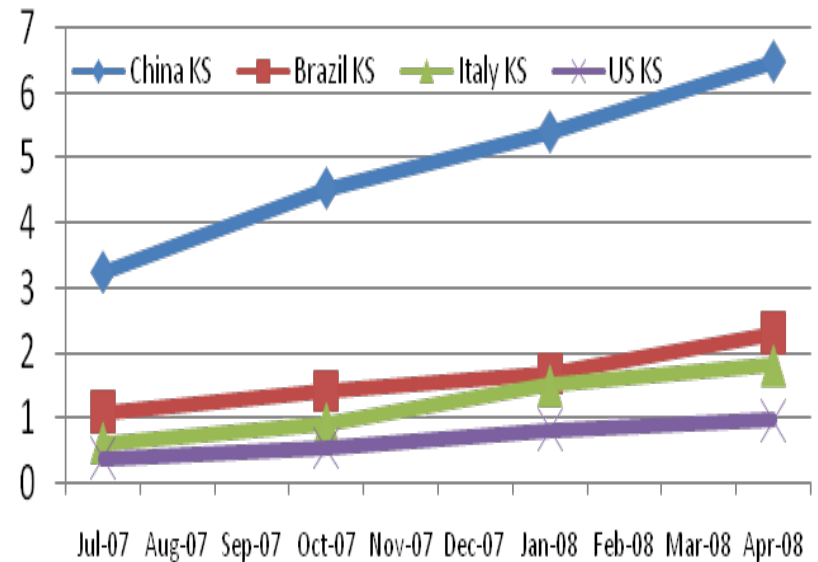
# Search Quality at Stake.

**61 countries have Q&A or advanced forums as top 10 most clicked destination**  
*(out of 115 countries with more than 1M session)*

% of First Result Page with >=1 Q&A Result from Yahoo or Baidu



% of Referral Traffic From 1<sup>st</sup> Page Sent to Yahoo / Baidu

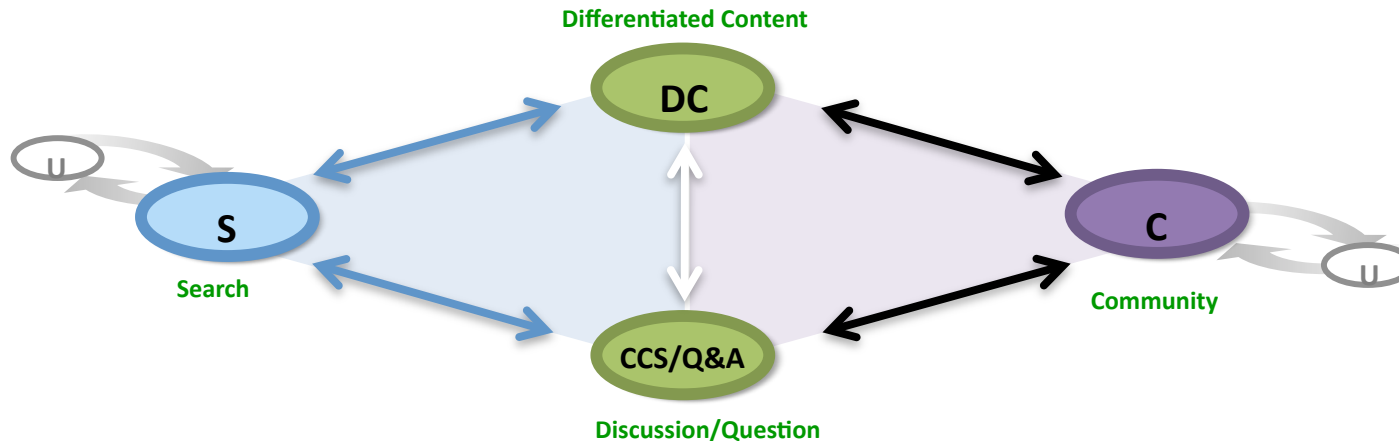


# SNS & Mobile Also Need Q&A

- Social Networks
  - Difficult to find user intent to match ads
  - Q&A is a perfect app to learn users' problems
- Mobile Search
  - Voice is the most convenient user interface
  - Succinct search result (or rich snippets) is desirable

# Confucius: Google Q&A

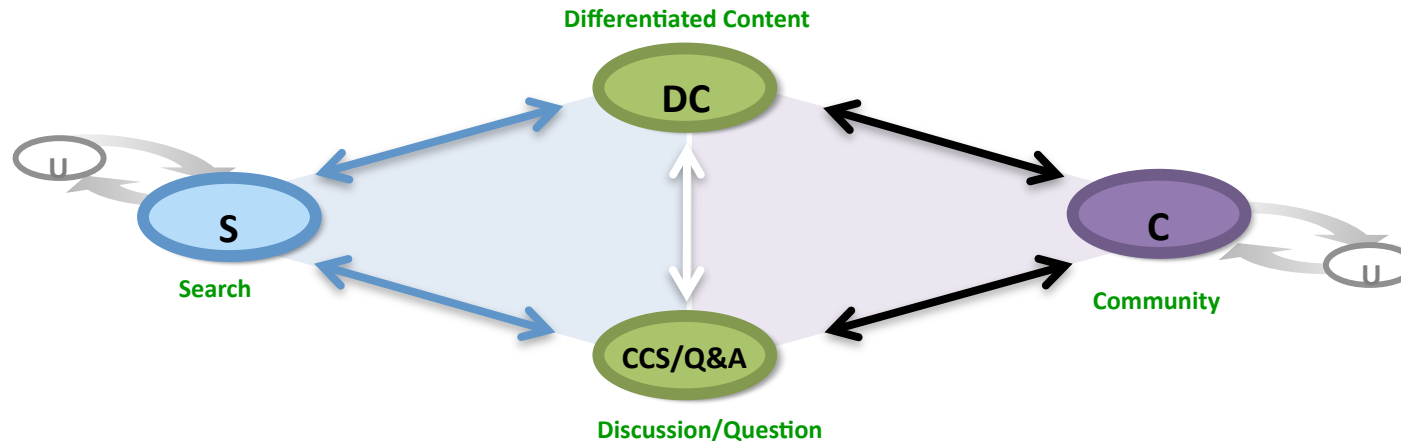
## Providing High-Quality Answers in a Timely Fashion



- Trigger a discussion/question session during search
- Provide labels to a post (semi-automatically)
- Given a post, find similar posts (automatically)
- Evaluate quality of a post, relevance and originality
- Evaluate user credentials in a topic sensitive way
- Route questions to experts
- Provide most relevant, high-quality content for Search to index
- Generate answers using NLP

# Confucius: Google Q&A

Providing High-Quality Answers in a Timely Fashion



- Trigger a discussion/question session during search
- Provide labels to a question (semi-automatically)
- Given a question, find similar questions (automatically)
- Evaluate quality of an answer, relevance and originality
- Evaluate user credentials in a topic sensitive way
- Route questions to experts
- Provide most relevant, high-quality content for Search to index
- Generate answers using NLP



首页 > 提问

提问的标题: iphone crack

详细描述:  
(选填)

Label suggestion using LDA algorithm.

悬赏问答分 10 你目前的问答分: 173

征答时限:  
(天) 10

添加标签:  
 电脑硬件  电脑软件  电脑基础  Windows  多  
 电脑游戏  网络游戏

• Real Time topic-to-topic (T2T) recommendation using LDA algorithm.

• Gives out related high quality links to previous questions before human answer appear.

已有的相关问答

- touch 2破解了吗 - 1个回答 60次浏览
- iphone 3g破解版如何上网 - 1个回答 940次浏览
- ipod touch2.2该不该破解? - 7个回答 69次浏览
- iphone视频存在哪个文件夹下? - 4个回答 74次浏览
- iPhone 3G2.2版本还用卡贴吗? - 1个回答 40次浏览
- iphone最新破解方法 - 1个回答 18次浏览
- iphone pc suite怎么用 - 3个回答 206次浏览
- 3G版iPhone是什么系统? 支持阅读PDF格式... - 1个回答 118次浏览

请选择1~5个与您的问题相关的标签 (需包含至少一个系统推荐的标签)

发表提问

请注意, 根据中国法律, 该服务会将有关您发帖内容、发帖时间以及您发帖时的IP地址、电子邮箱地址等记录保留至少 60 天, 并且只要接到合法请求, 即会将这类信息提供给政府机构。点击“发表提问”表示您接受服务

# Collaborative Filtering

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



Predict further *membership*

Labels/Qs

		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					

Questions

# FIM-based Recommendation



To grow the base, we need association rules

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

# Distributed Latent Dirichlet Allocation (LDA)

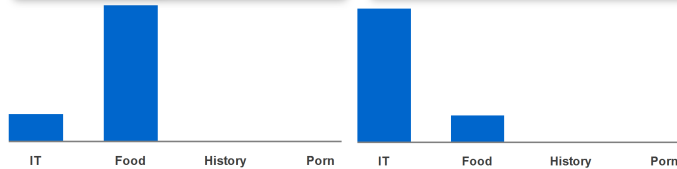
- Search
  - Construct a latent layer for better for semantic matching
- Example:
  - iPhone crack
  - Apple pie

Documents

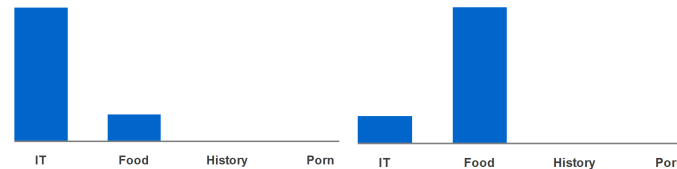
1 recipe pastry for a 9 inch double crust  
9 apples, 2/1 cup, brown sugar

How to install apps on Apple mobile phones?

Topic Distribution



Topic Distribution



User queries

iPhone crack

Apple pie

2010-12-13

WISE Keynote

Users/Music/Ads/Question

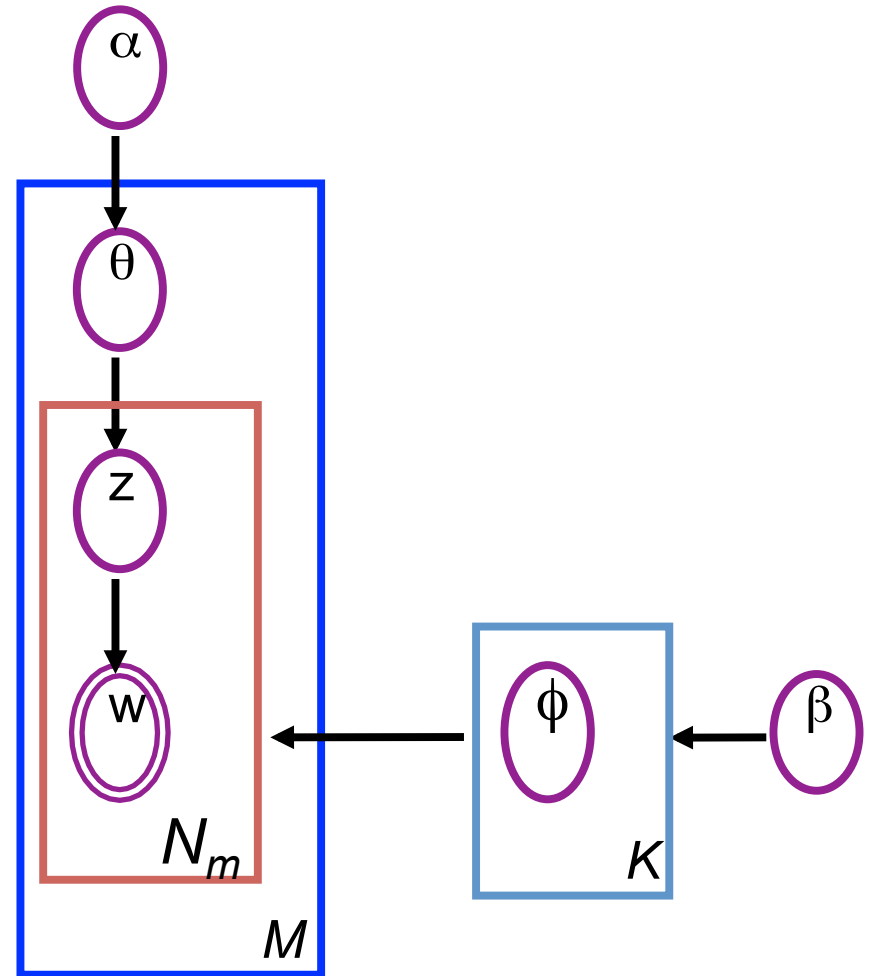
?	?	1	3	1	?	?	?	?	?	?
?	2	?	1	2	?	1	?	3	?	1
?	?	?	?	?	1		5			1
	5		3		1	1				
		1								
						1	4			
			2					1		
1	2									
	1								5	
1										1
	1	4	1	3	6					

Users/Music/Ads/Answers

- Other Collaborative Filtering Apps
  - Recommend Users → Users
  - Recommend Music → Users
  - Recommend Ads → Users
  - Recommend Answers → Q
- Predict the ? In the light-blue cells

# Latent Dirichlet Allocation [D. Blei, M. Jordan 04]

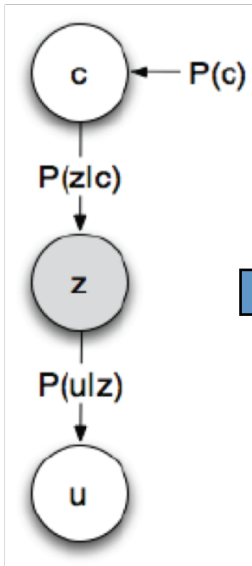
- $\alpha$ : uniform Dirichlet  $\phi$  prior for per document  $d$  topic distribution (corpus level parameter)
- $\beta$ : uniform Dirichlet  $\phi$  prior for per topic  $z$  word distribution (corpus level parameter)
- $\theta_d$  is the topic distribution of document  $d$  (document level)
- $z_{dj}$  the topic if the  $j^{\text{th}}$  word in  $d$ ,  $w_{dj}$  the specific word (word level)



# Combinational Collaborative Filtering Model (CCF)

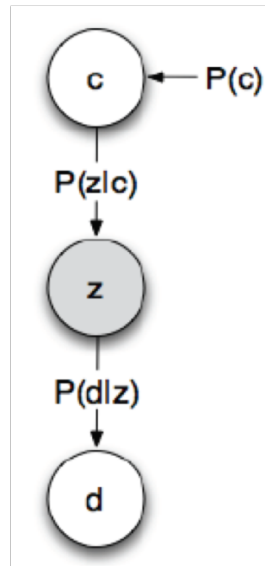
[W.-Y. Chen, et al, KDD2008]

Communities

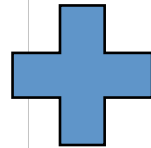


users

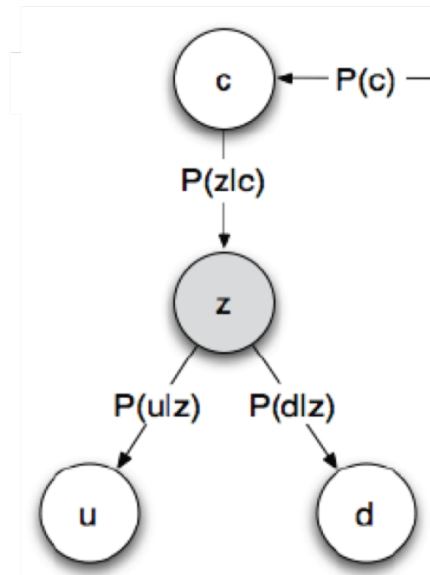
Communities



descriptions

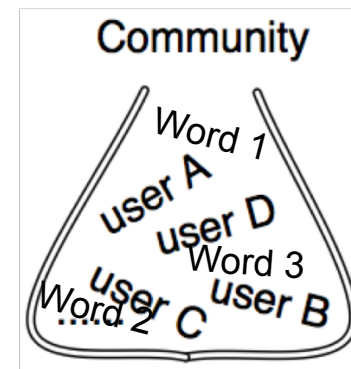


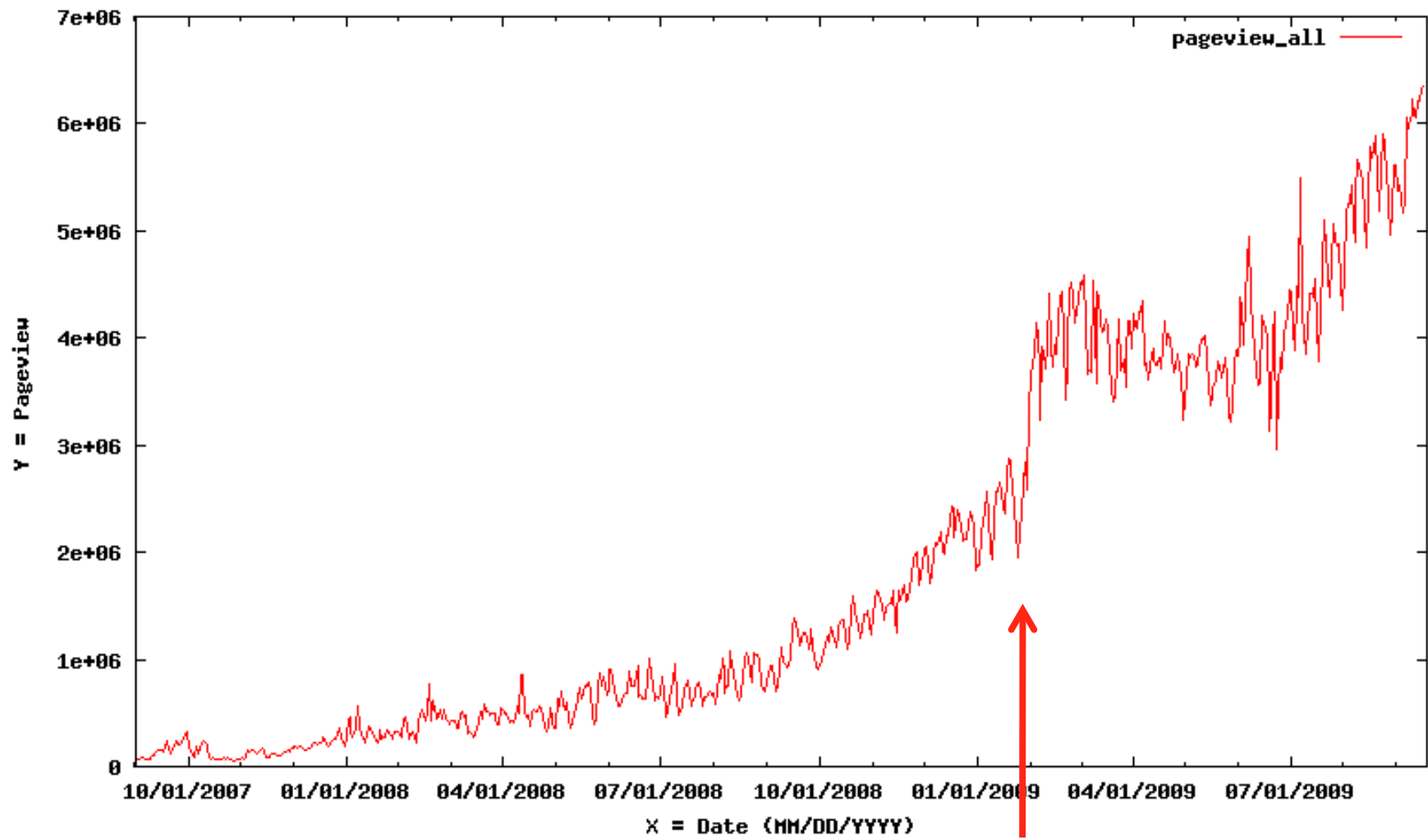
Communities



users

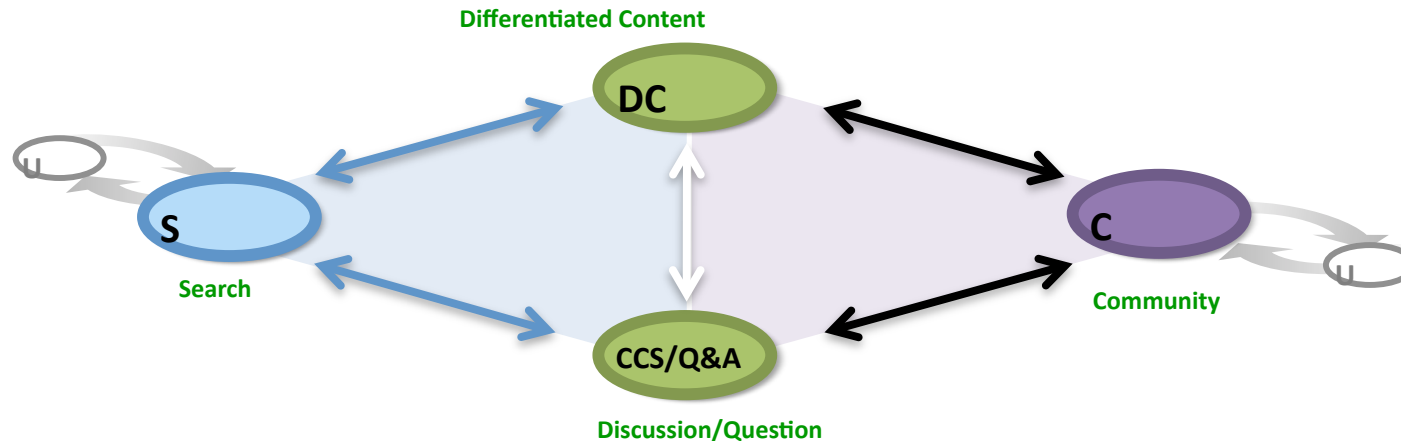
descriptions





Generated on 2009-09-19

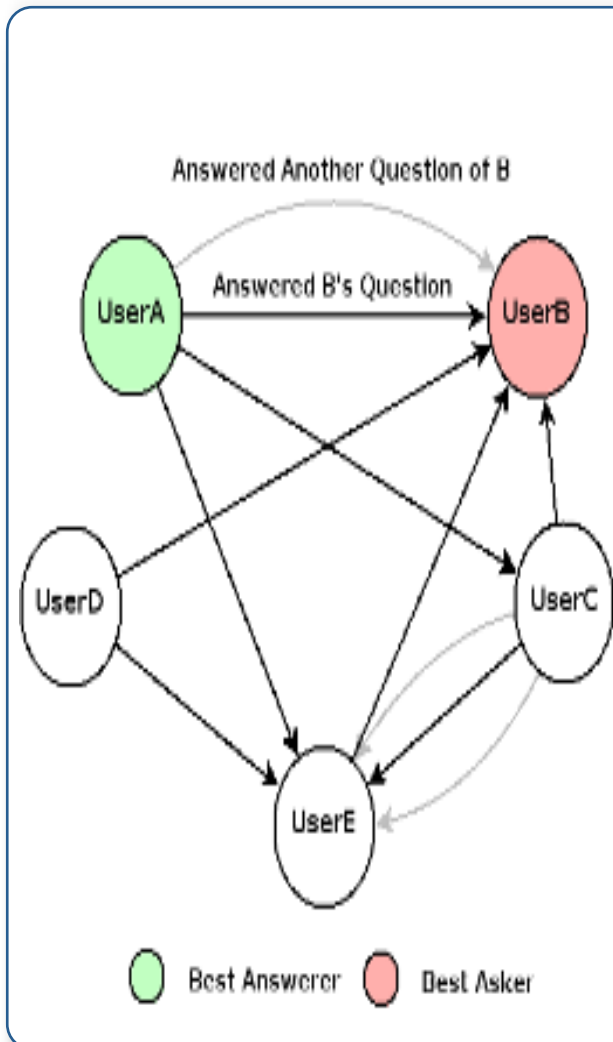
# Confucius: Google Q&A



- Trigger a discussion/question session during search
- Provide labels to a post (semi-automatically)
- Given a post, find similar posts (automatically)
- Evaluate quality of a post, relevance and originality
- Evaluate user credentials in a topic sensitive way
- Route questions to experts
- Provide most relevant, high-quality content for Search to index
- NLQA



# UserRank



- Rank users by quantity (**number of links**) and quality (**weights on links**) of contributions

Quality include:

- **Relevance.** Is an answer relevant to the Q? Measured by KL divergence between *latent-topic vectors* of A and Q
- **Coverage.** Compared among different answers
- **Originality.** Detect potential plagiarism and spam
- **Promptness.** Time between Q and A posting time

# Outline

- Search + Social Synergy
- Social → Search
- Search → Social
- Scalability

# Outline

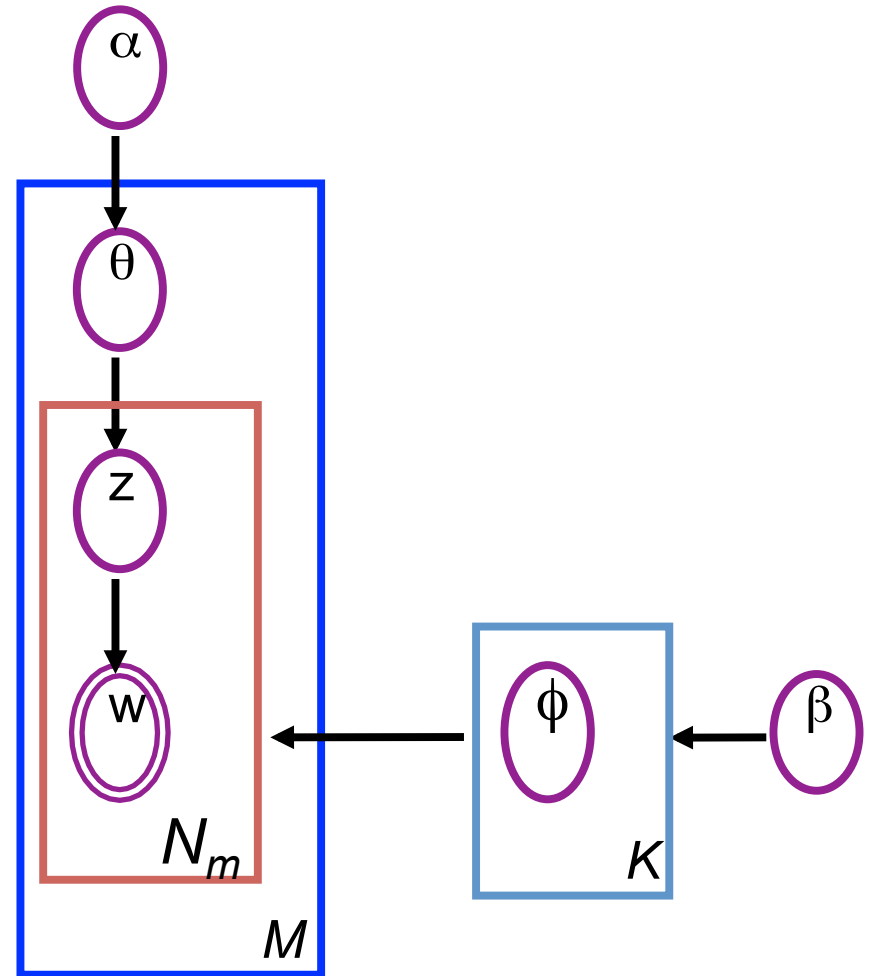
- Search + Social Synergy
- Social → Search
- • Search → Social
- Scalability

# Social?

- Connecting to friends
- Knowing what friends are up to
- Connecting to strangers
  - Dating, Gaming
  - Shopping
- Making recommendations based on activities

# User Latent Model

- $\alpha$ : uniform Dirichlet  $\phi$  prior for per user  $u$  interest distribution (population level parameter)
- $\beta$ : uniform Dirichlet  $\phi$  prior for per interest  $z$  activity distribution (population level parameter)
- $\theta_d$  is the interest distribution of user  $u$  (user level)
- $z_{uj}$  the interest of the  $j^{\text{th}}$  activity in  $u$ ,  $w_{uj}$  the specific activity (activity level)



?	?	1	3	1	?	?	?	?	?
?	2	?	1	2	?	1	?	3	?
?	?	?	?	?	1		5		
	5		3		1	1			
		1							
					1	4			
							1		
			2						
1	2							5	
	1								1
1									
	1	4	1	3	6				

?	1	3	1	?	?	?	?	?	?
2	?	1	2	?	1	?	3	?	1
?	?	?	?	?	1		5		
5		3			1	1			
	1								
					1	4			
							1		
			2						
2								5	
1									1
	1	4	1	3	6				

?	?	1	3	1	?	?	?	?	?
?	2	?	1	2	?	1	?	3	?
?	?	?	?	?	1		5		
?	5		3		1	1			
		1							
					1	4			
							1		
			2						
1	2							5	
	1								
1									
	1	4	1	3	6				

?	?	1	3	1	?	?	?	?	?
?	2	?	1	2	?	1	?	3	?
?	?	?	?	?	1		5		
?	?	?	?	?	1		5		
		1							
					1	4			
							1		
			2						
1	2							5	
	1								1
	1	4	1	3	6				

?	?	1	3	1	?	?	?	?	?
?	2	?	1	2	?	1	?	3	?
?	?	?	?	?	1		5		
?	?	?	?	?	1		5		
		1							
					1	4			
							1		
			2						
1	2							5	
	1								
1									
	1	4	1	3	6				

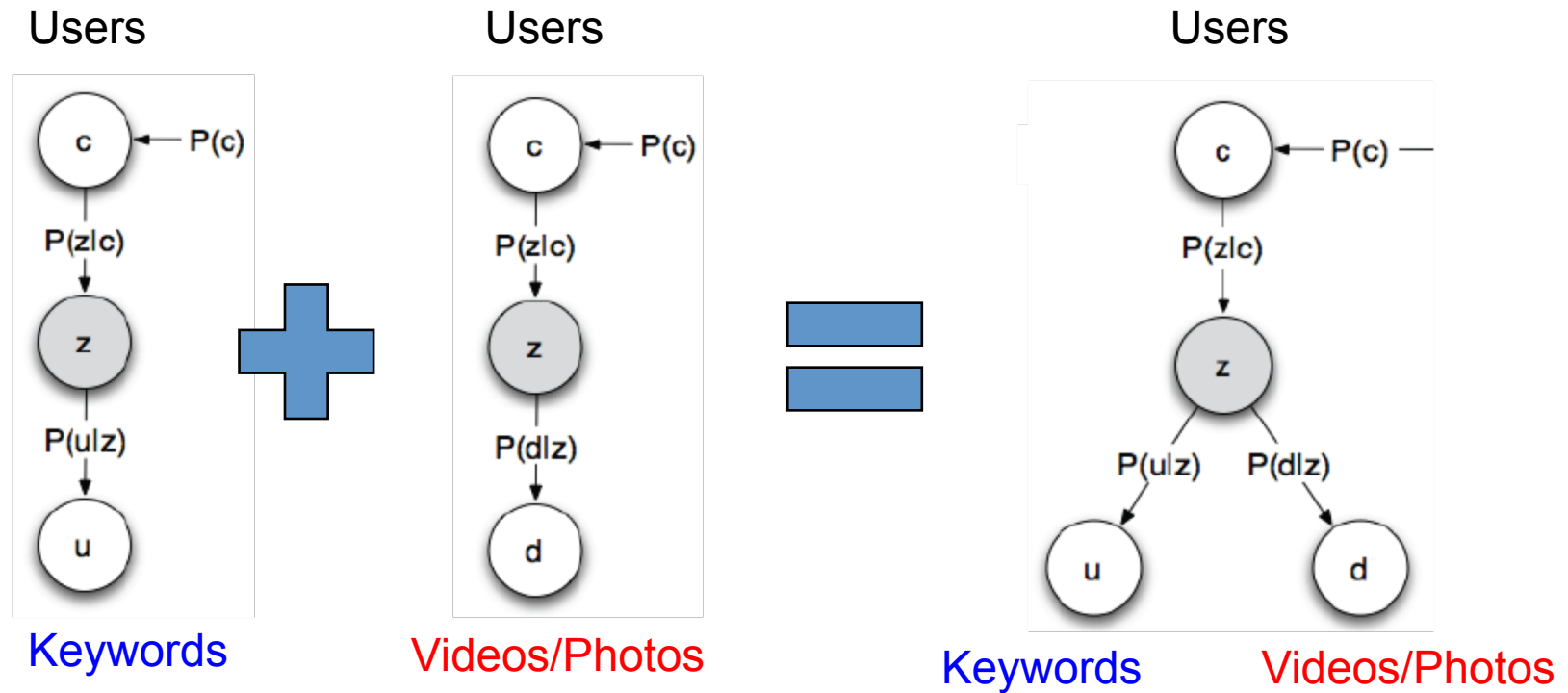
?	?	1	3	1	?	?	?	?	?
?	2	?	1	2	?	1	?	3	?
?	?	?	?	?	1		5		
	5		3		1	1			
		1							
					1	4			
							1		
			2						
1	2							5	
	1								1
1									
	1	4	1	3	6				

?	?	1	3	1	?	?	?	?	?
?	2	?	1	2	?	1	?	3	?
?	?	?	?	?	1		5		
?	?	?	?	?	1		5		
		1							
					1	4			
							1		
			2						
1	2							5	
	1								1
1									
	1	4	1	3	6				

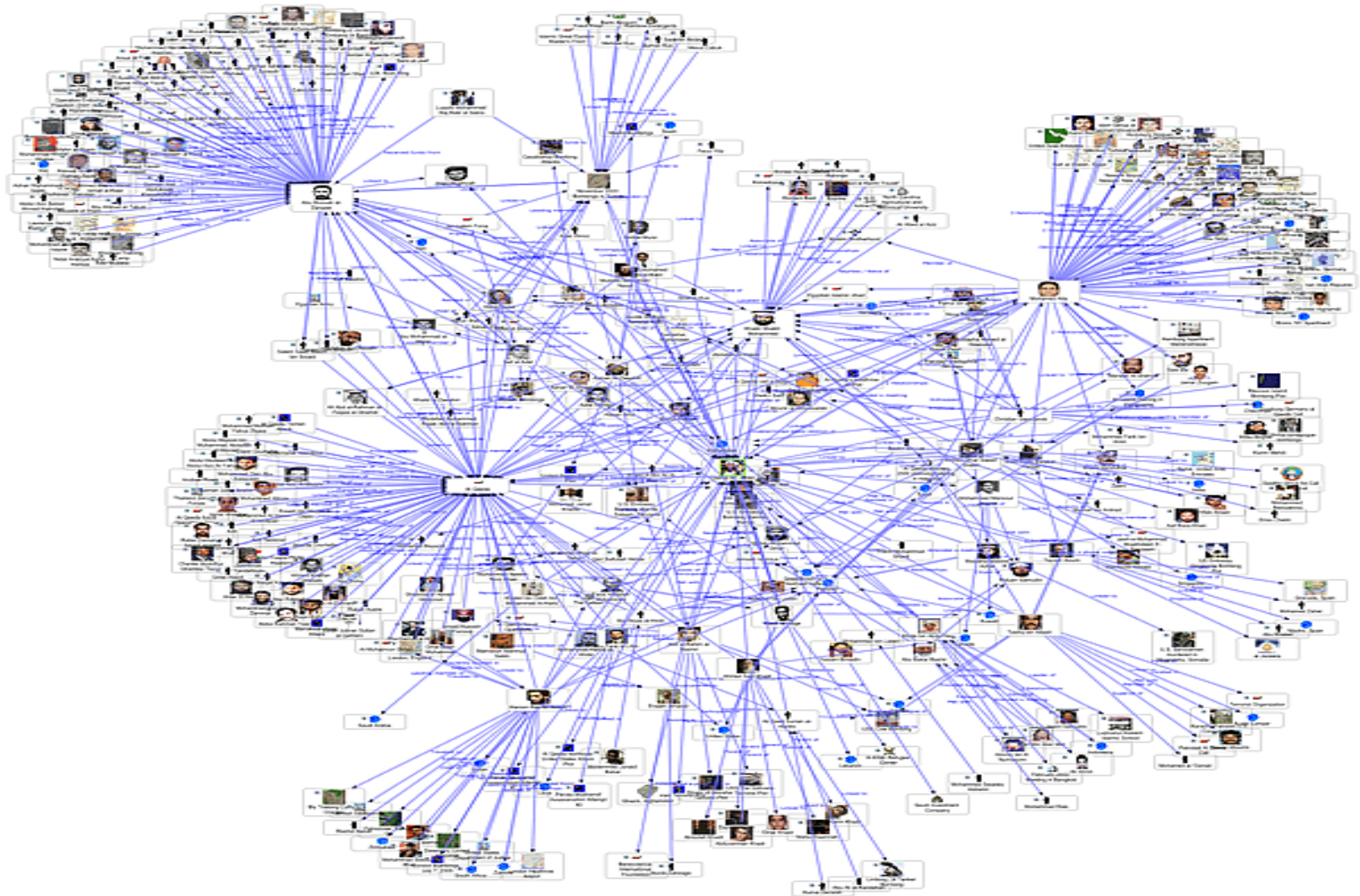
2010-12-13

WISE Keynote

# Combinational Collaborative Filtering Model (CCF)



# Interest Networks





# Outline

- Search + Social Synergy
- Social → Search
  - Mobilize users to improve search quality
  - Google *Q&A*, Facebook *Like*
- Search → Social
  - Use query log to help social
    - Activities → Interests → Social
    - Groupcom
- • Scalability

# Prefixes

SI prefix	Name	Power of 10 or 2	
k kilo	thousand	$10^3$ $2^{10}$	
M mega	million	$10^6$ $2^{20}$	
G giga	billion	$10^9$ $2^{30}$	
T tera	trillion	$10^{12}$ $2^{40}$	
P peta	quadrillion	$10^{15}$ $2^{50}$	
E exa	quintillion	$10^{18}$ $2^{60}$	
Z zetta	sextillion	$10^{21}$ $2^{70}$	
Y yotta	septillion	$10^{24}$ $2^{80}$	

# Prefixes

SI prefix	Name	Power of 10 or 2	
k kilo	thousand	$10^3$ $2^{10}$	
M mega	million	$10^6$ $2^{20}$	
G giga	billion	$10^9$ $2^{30}$	
T tera	trillion	$10^{12}$ $2^{40}$	
P peta	quadrillion	$10^{15}$ $2^{50}$	
E exa	quintillion	$10^{18}$ $2^{60}$	
Z zetta	sextillion	$10^{21}$ $2^{70}$	
Y yotta	septillion	$10^{24}$ $2^{80}$	

# Prefixes

SI prefix	Name	Power of 10 or 2	
k kilo	thousand	$10^3$ $2^{10}$	
M mega	million	$10^6$ $2^{20}$	
G giga	billion	$10^9$ $2^{30}$	
T tera	trillion	$10^{12}$ $2^{40}$	
P peta	quadrillion	$10^{15}$ $2^{50}$	
E exa	quintillion	$10^{18}$ $2^{60}$	
Z zetta	sextillion	$10^{21}$ $2^{70}$	
Y yotta	septillion	$10^{24}$ $2^{80}$	

# More Data vs. Better Algorithms

Banko & Brill, 2001

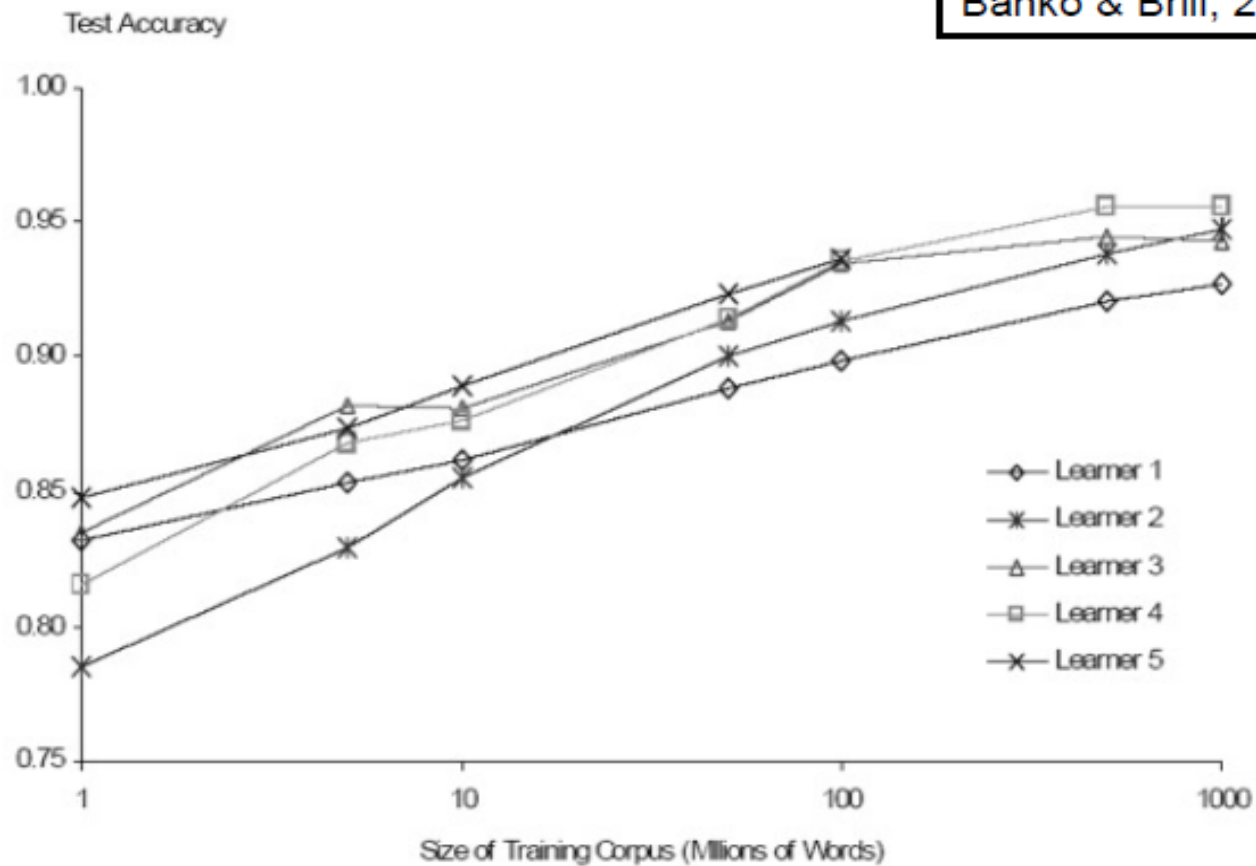


Figure 2. Learning Curves for Confusable Disambiguation

# More Data vs. Better Algorithms

Banko & Brill, 2001

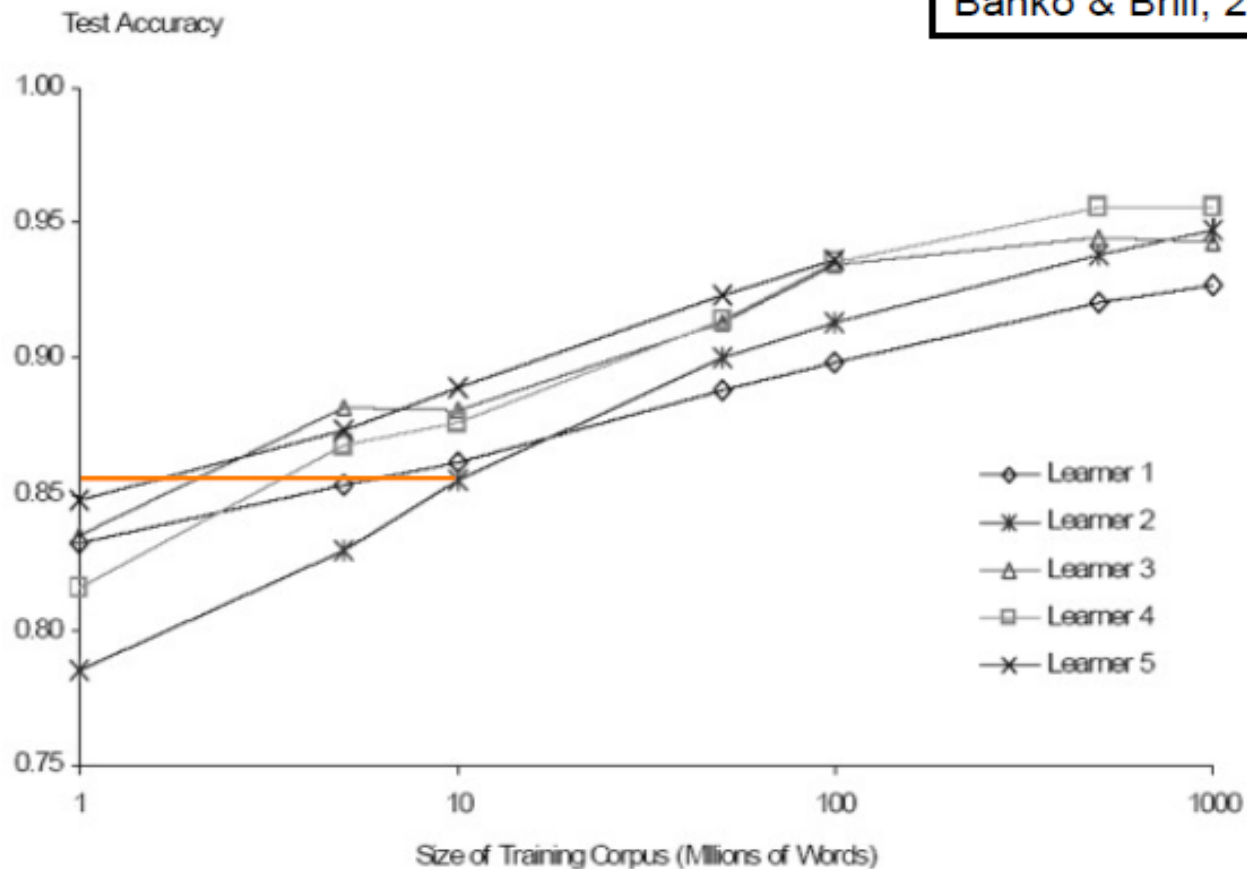


Figure 2. Learning Curves for Confusable Disambiguation

# More Data vs. Better Algorithms

Banko & Brill, 2001

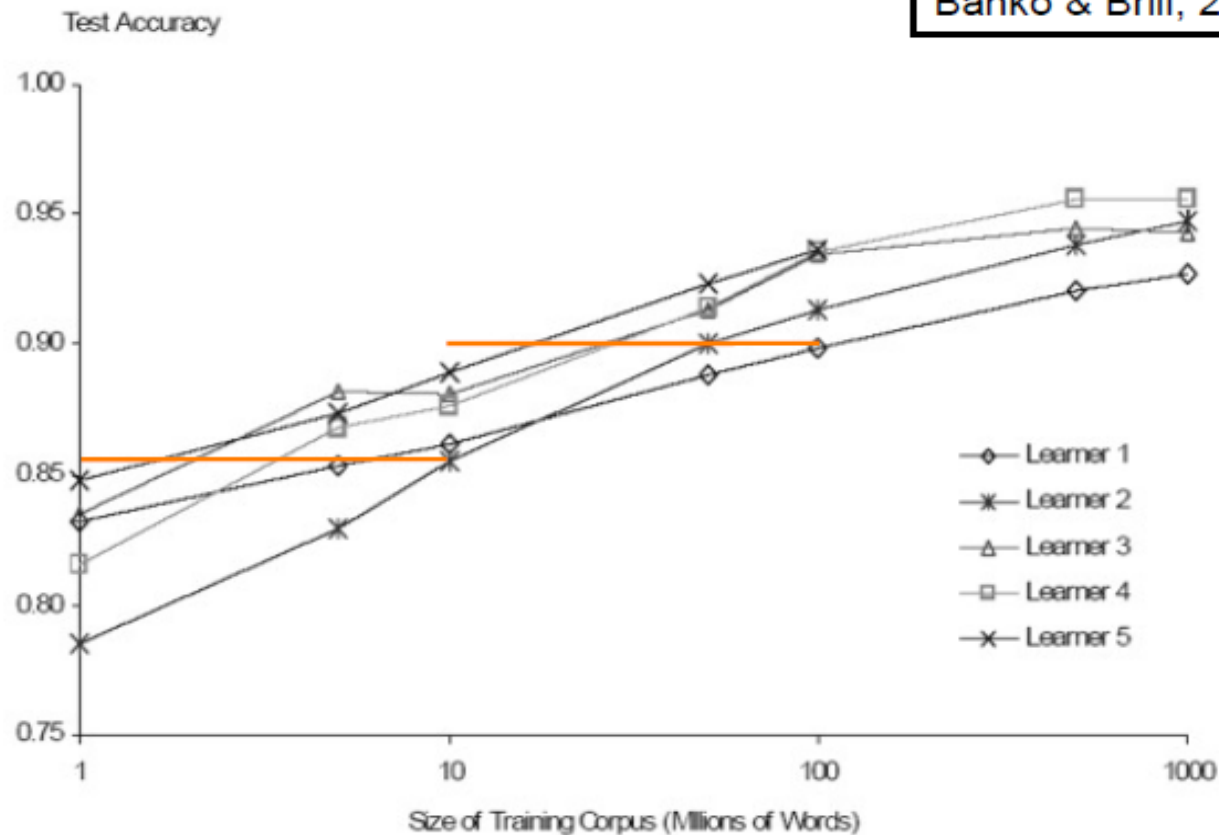


Figure 2. Learning Curves for Confusable Disambiguation

# More Data vs. Better Algorithms

Banko & Brill, 2001

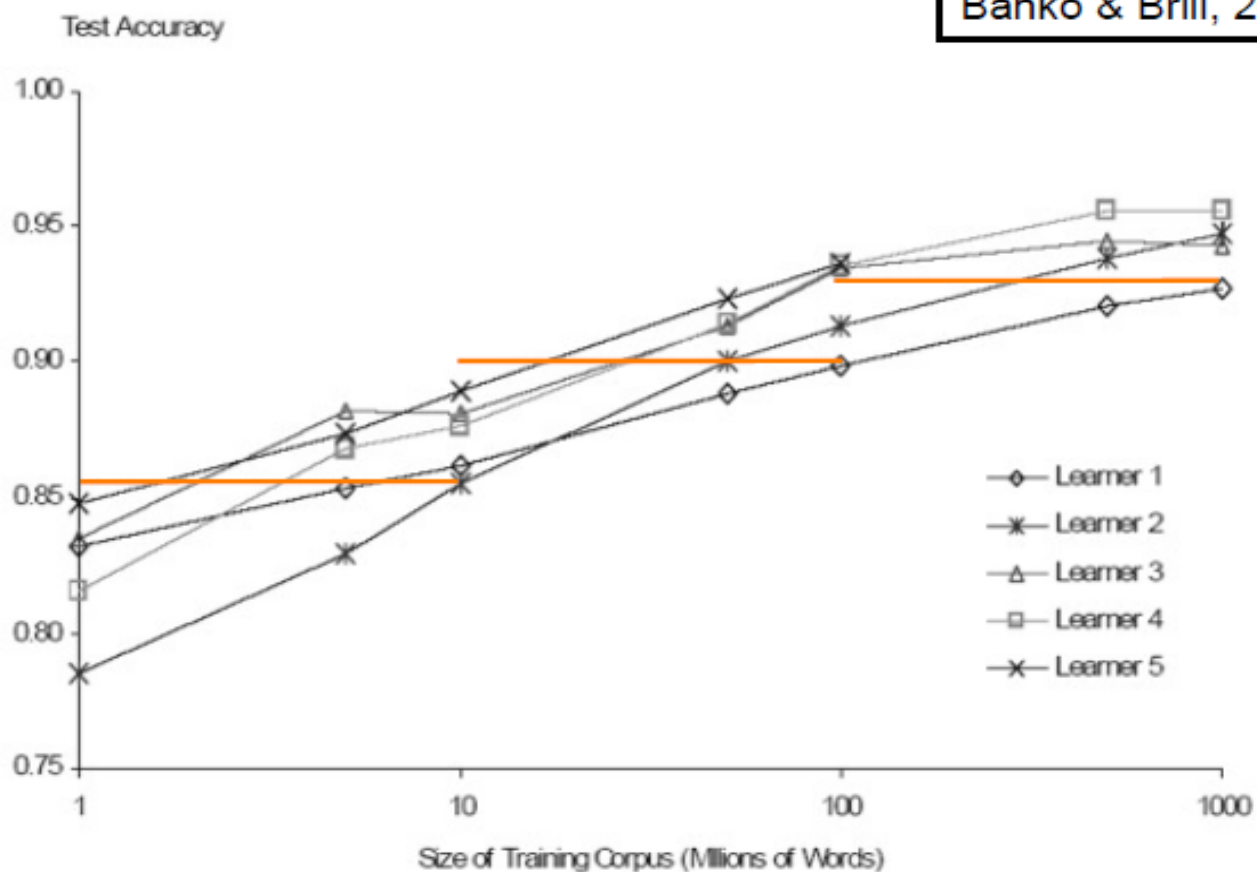
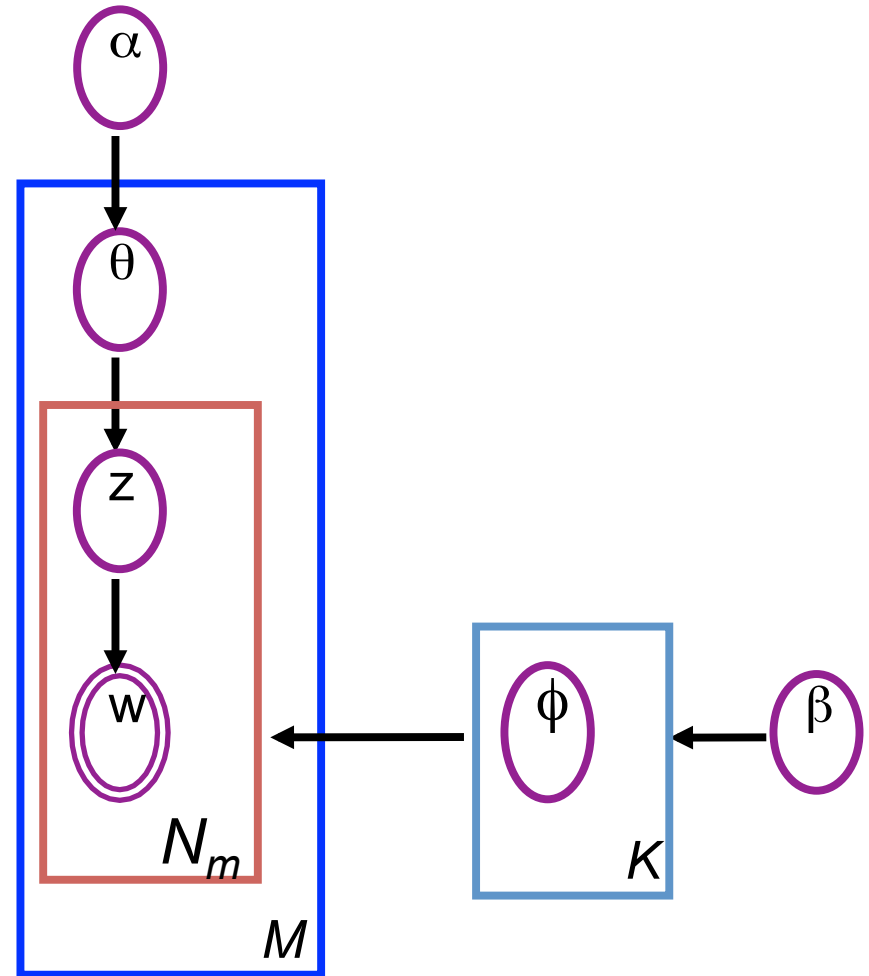


Figure 2. Learning Curves for Confusable Disambiguation



# User Latent Model

- $\alpha$ : uniform Dirichlet  $\phi$  prior for per user  $u$  interest distribution (population level parameter)
- $\beta$ : uniform Dirichlet  $\phi$  prior for per interest  $z$  activity distribution (population level parameter)
- $\theta_d$  is the interest distribution of user  $u$  (user level)
- $z_{uj}$  the interest of the  $j^{\text{th}}$  activity in  $u$ ,  $w_{uj}$  the specific activity (activity level)



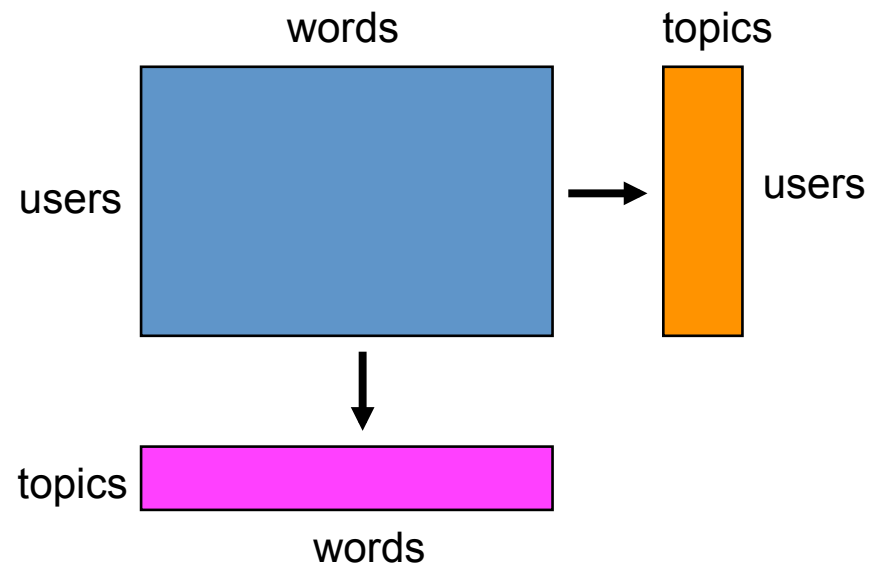
# LDA Gibbs Sampling: Inputs & Outputs

## Inputs:

1. training data: users as bags of words
2. parameter: the number of topics

## Outputs:

1. model parameters: a co-occurrence matrix of topics and words.
2. by-product: a co-occurrence matrix of topics and users.



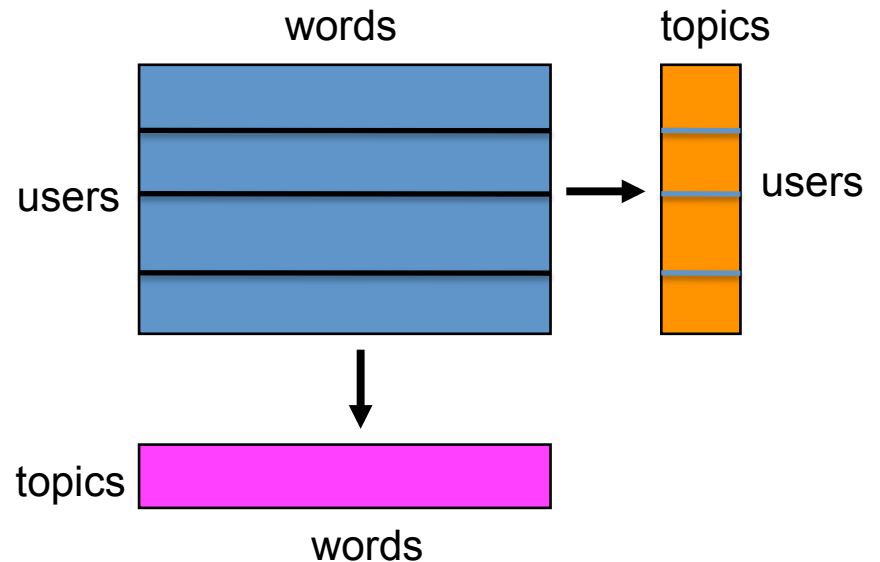
# Parallel Gibbs Sampling

## Inputs:

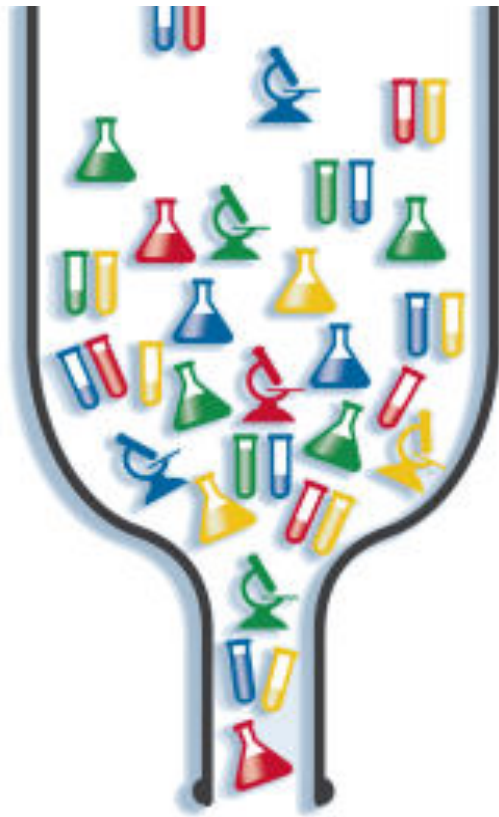
1. training data: users as bags of words
2. parameter: the number of topics

## Outputs:

1. model parameters: a co-occurrence matrix of topics and words.
2. by-product: a co-occurrence matrix of topics and users.



# Observations



Master Node

- Bottleneck:  
Communication
- Amdahl's law caps speedup
- Words in a bag have no order
- Words on a computer node can be reordered

# Example Bags / Node A

- Bag #1 w1, w2, w3, w1, w2, w3, w1, w2, w3
  - Bag #2 w1, w2, w1, w2, w1, w2, w1, w2
  - Bag #3 w3, w1, w3, w1, w3, w1, w3, w1
- 
- Bundle #1 w1, w1, w1, w1, w1, w1, ...
  - Bundle #2 w2, w2, w2, ... ,
  - Bundle #3 w3, w3, w3, ...

# Two Nodes

Node A	Node B
W1	W2
W2	W3
W3	W1

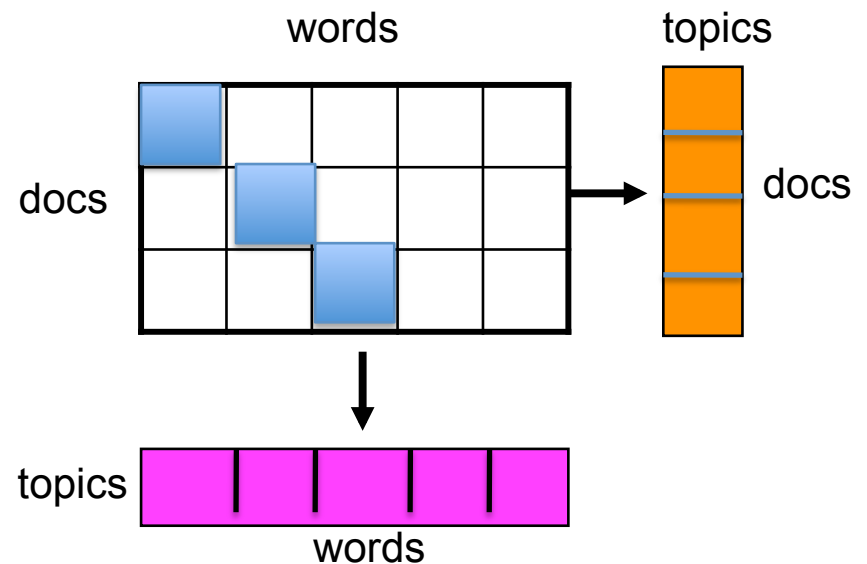
# Parallel Gibbs Sampling

## Inputs:

1. training data: documents as bags of words
2. parameter: the number of topics

## Outputs:

1. model parameters: a co-occurrence matrix of topics and words.
2. by-product: a co-occurrence matrix of topics and documents.



# PLDA -- enhanced parallel LDA

- Take advantage of bag of words modeling: each Pw machine processes vocabulary in a word order
- Pipelining: fetching the updated topic distribution matrix while doing Gibbs sampling

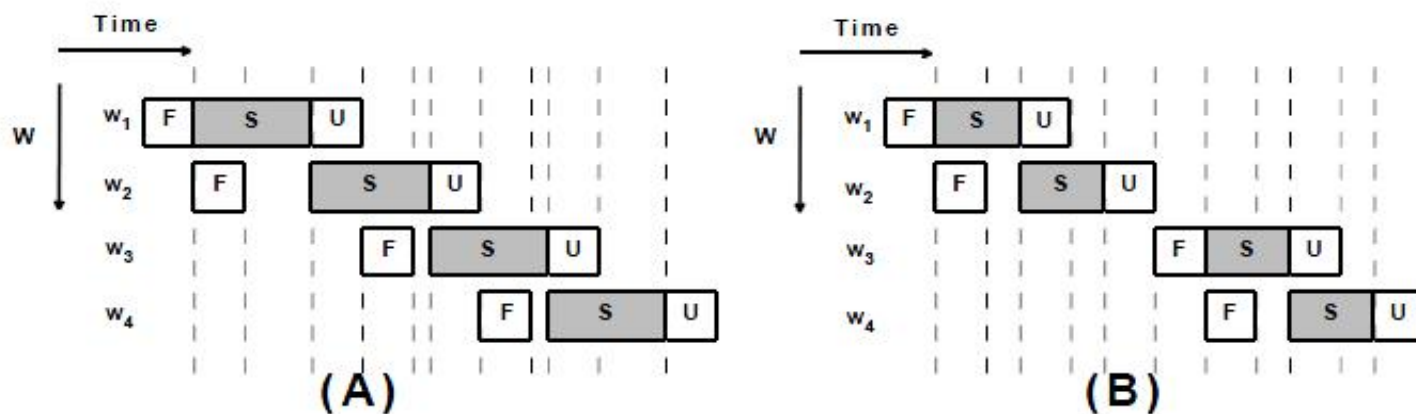
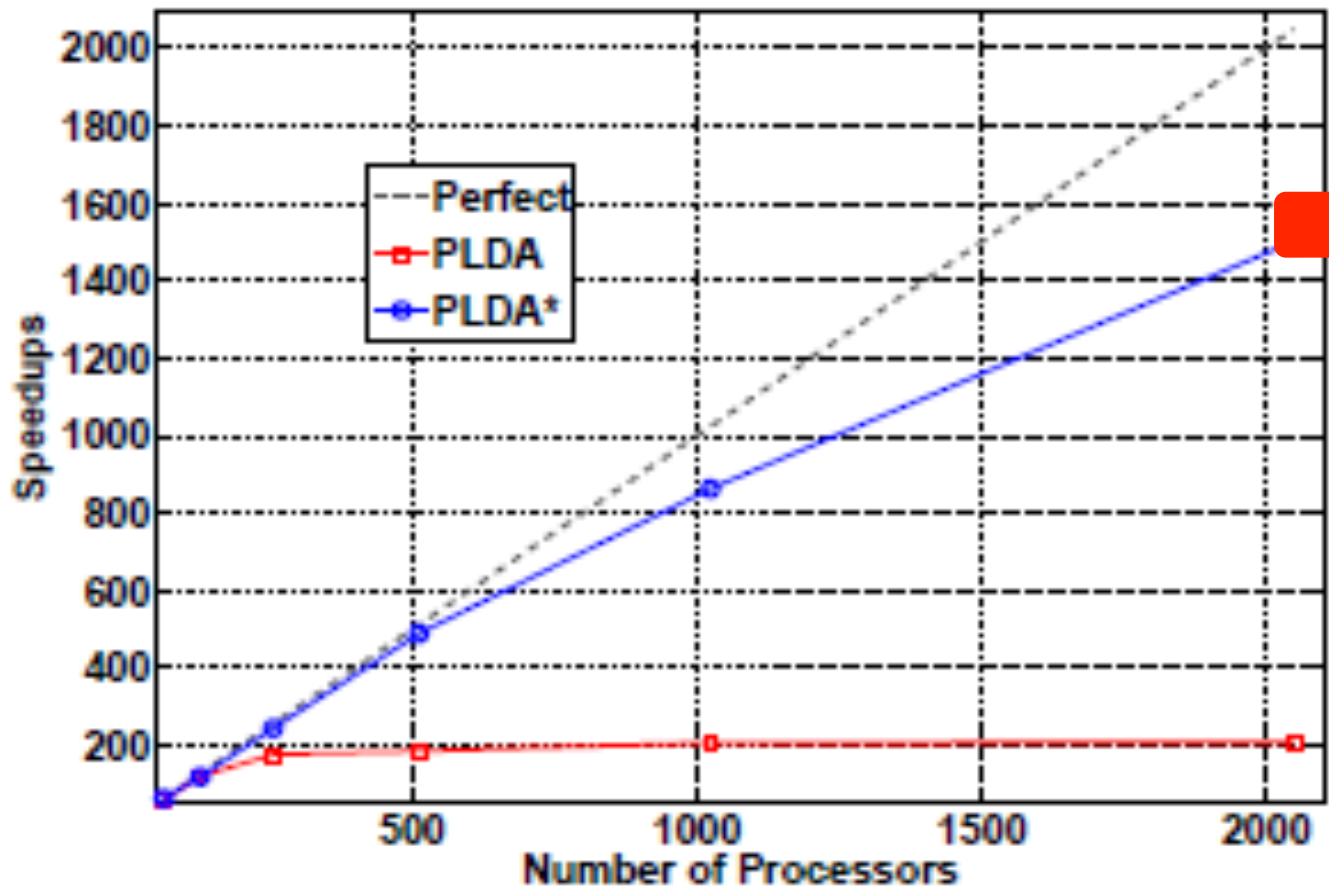


Fig. 4: Pipeline-based Gibbs Sampling in PLDA\*. (A):  $t_s \geq t_f + t_u$ . (B):  $t_s < t_f + t_u$ .



# Speedup

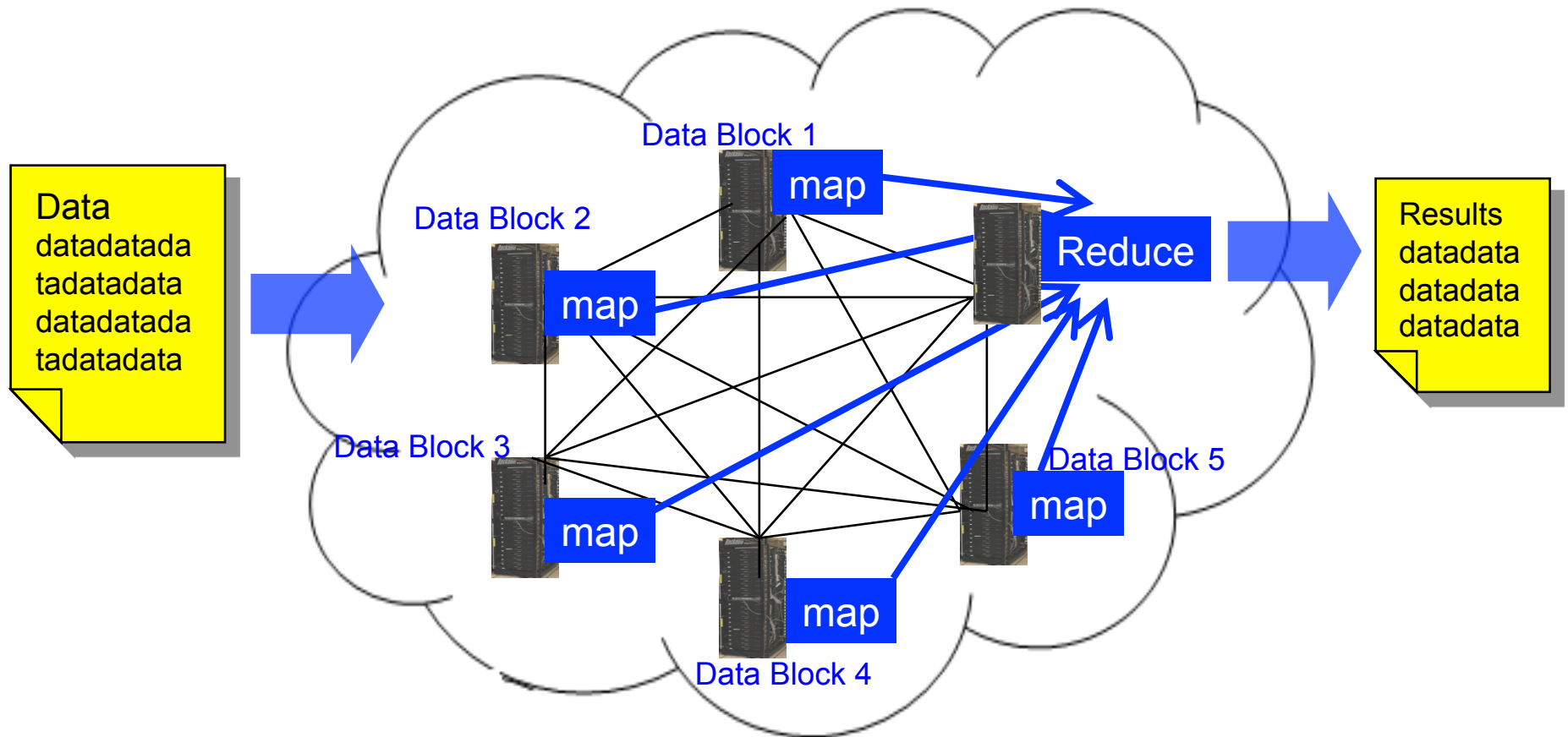
1,500x using 2,000 machines



# Lessons Learned

- Bottleneck Matters
- Inter-iteration Matters

# MapReduce

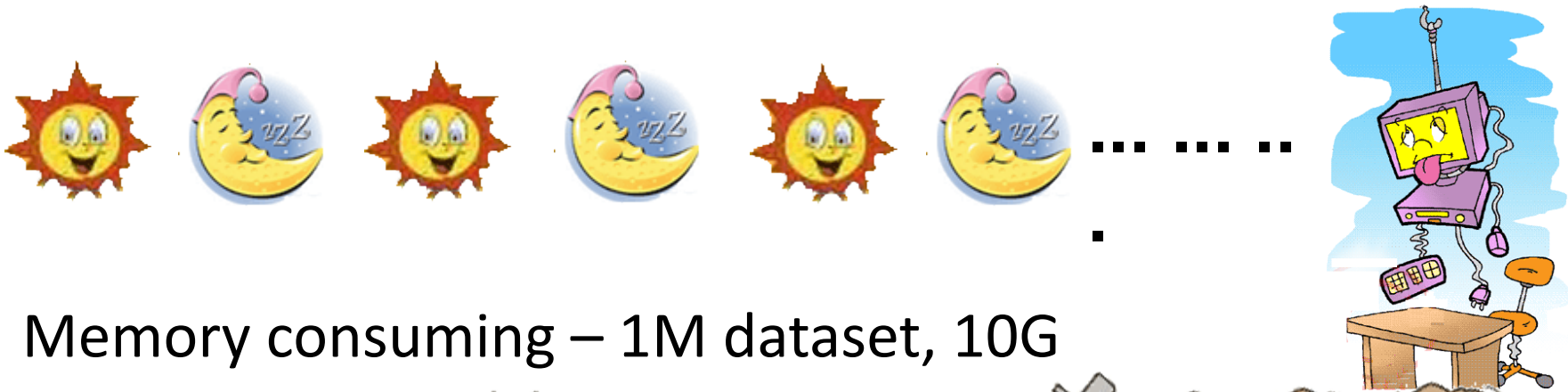


# Parallel Programming Models

	MapReduce	Project +	MPI
GFS/IO and task rescheduling overhead between iterations	Yes	No +1	No +1
Flexibility of computation model	AllReduce only +0.5	? +1	Flexible +1
Efficient AllReduce	Yes +1	Yes +1	Yes +1
Recover from faults between iterations	Yes +1	Yes +1	Apps
Recover from faults within each iteration	Yes +1	Yes +1	Apps
Final Score for scalable machine learning	3.5	5	4

# SVM Bottlenecks

Time consuming – 1M dataset, 8 days



Memory consuming – 1M dataset, 10G



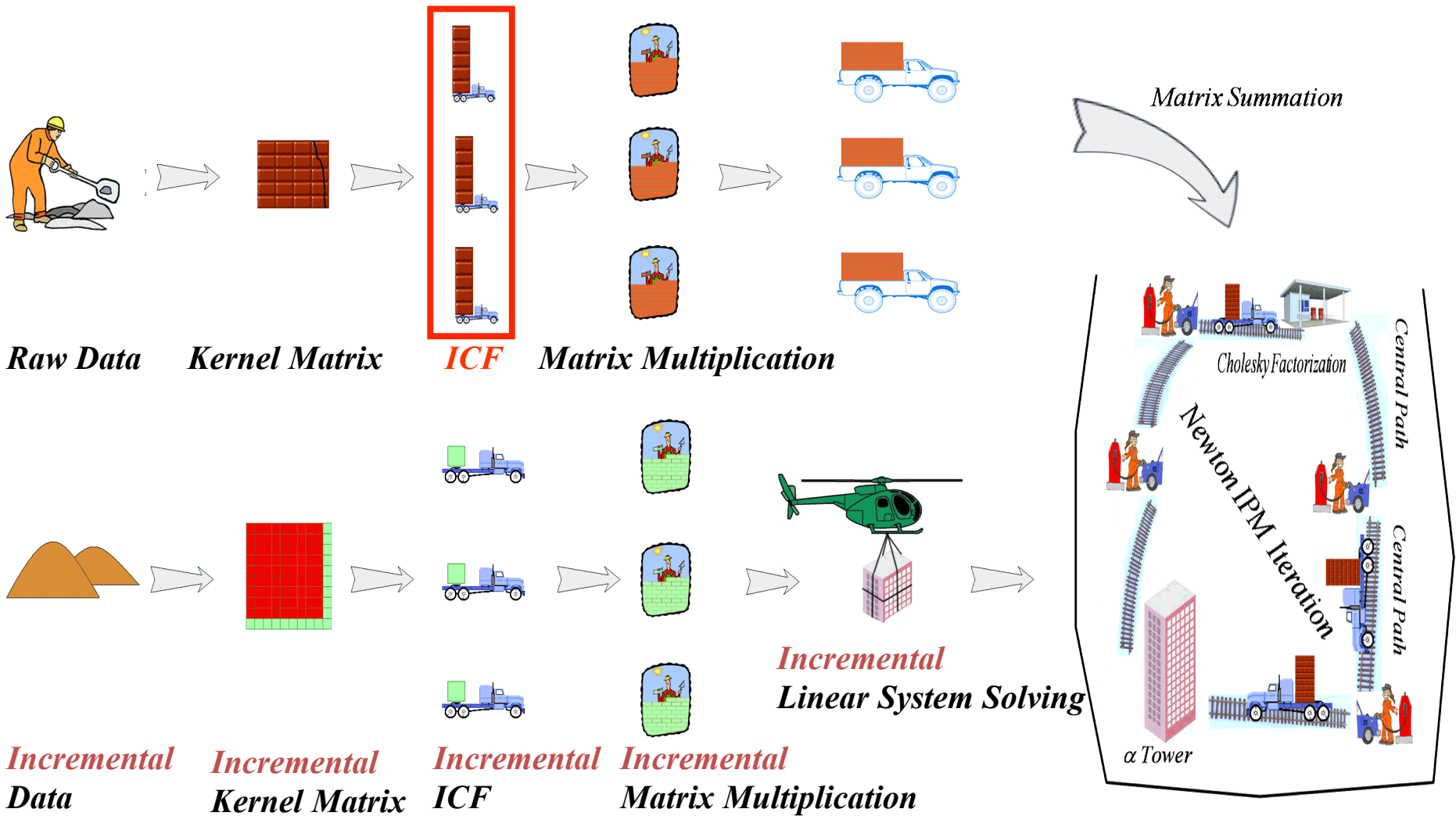
# Matrix Factorization Alternatives

Factorization	Cost
QR	$O(\frac{4}{3}n^3)$
LU	$O(\frac{2}{3}n^3)$
Cholesky	$O(\frac{1}{3}n^3 + 2n^2)$
LDL <sup>T</sup>	$O(\frac{1}{3}n^3)$
Incomplete Cholesky	$O(p^2n)$
Kronecker	$O(2n^2)$

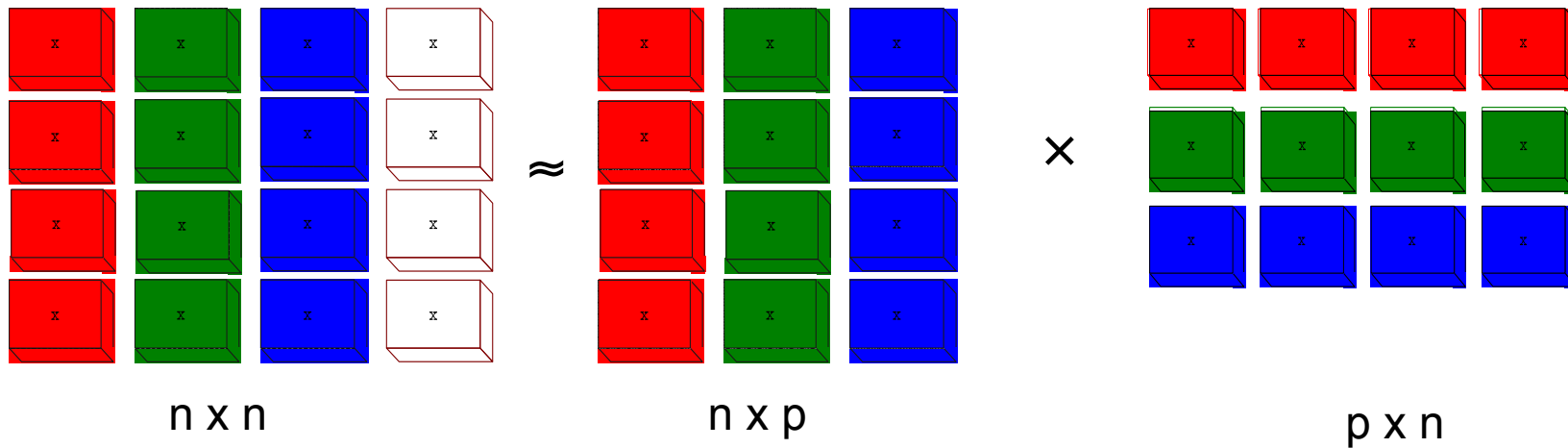
exact ←

→ approximate

# Parallelizing SVM [E. Chang, et al, NIPS 07]

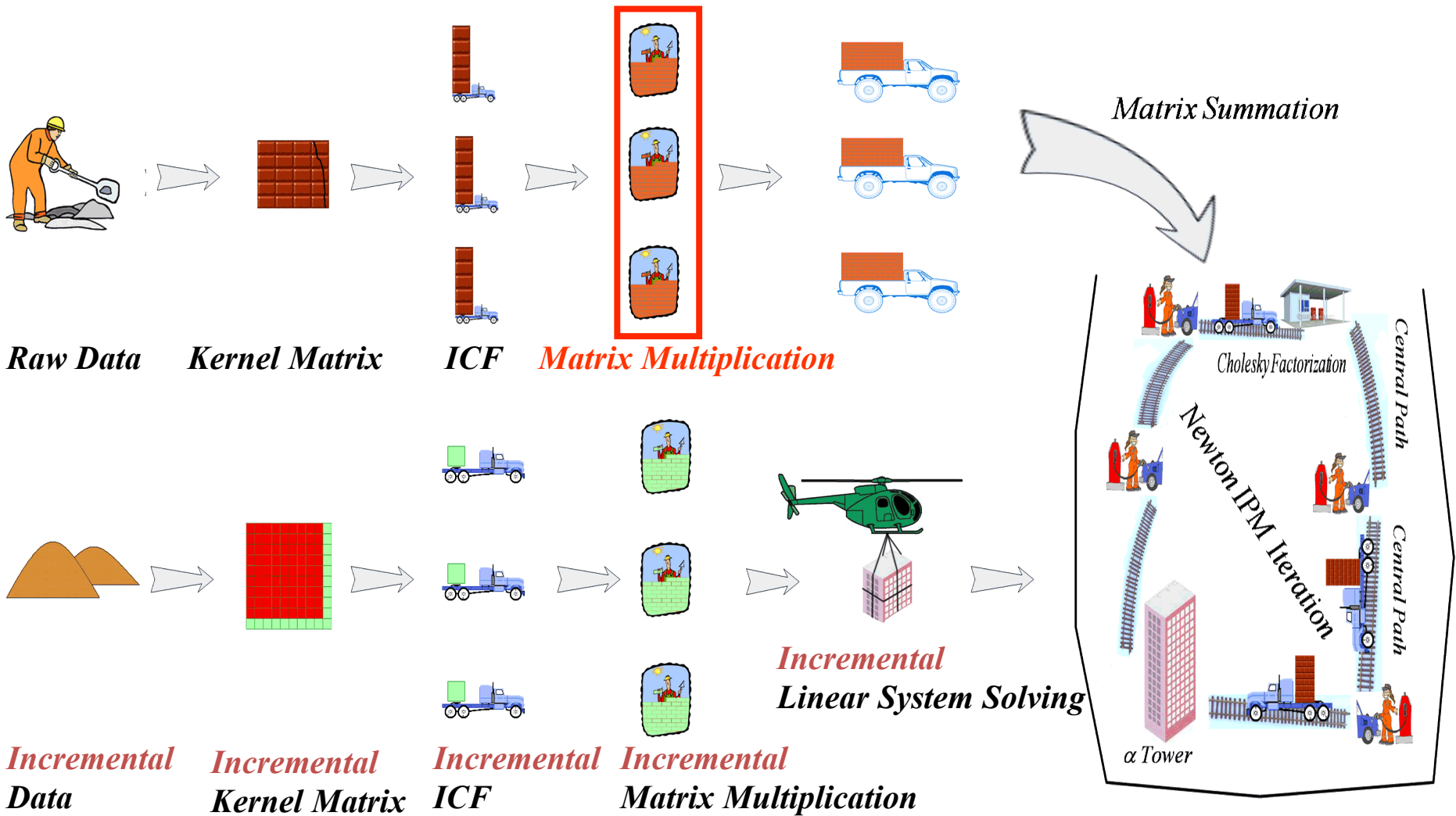


# Incomplete Cholesky Factorization (ICF)

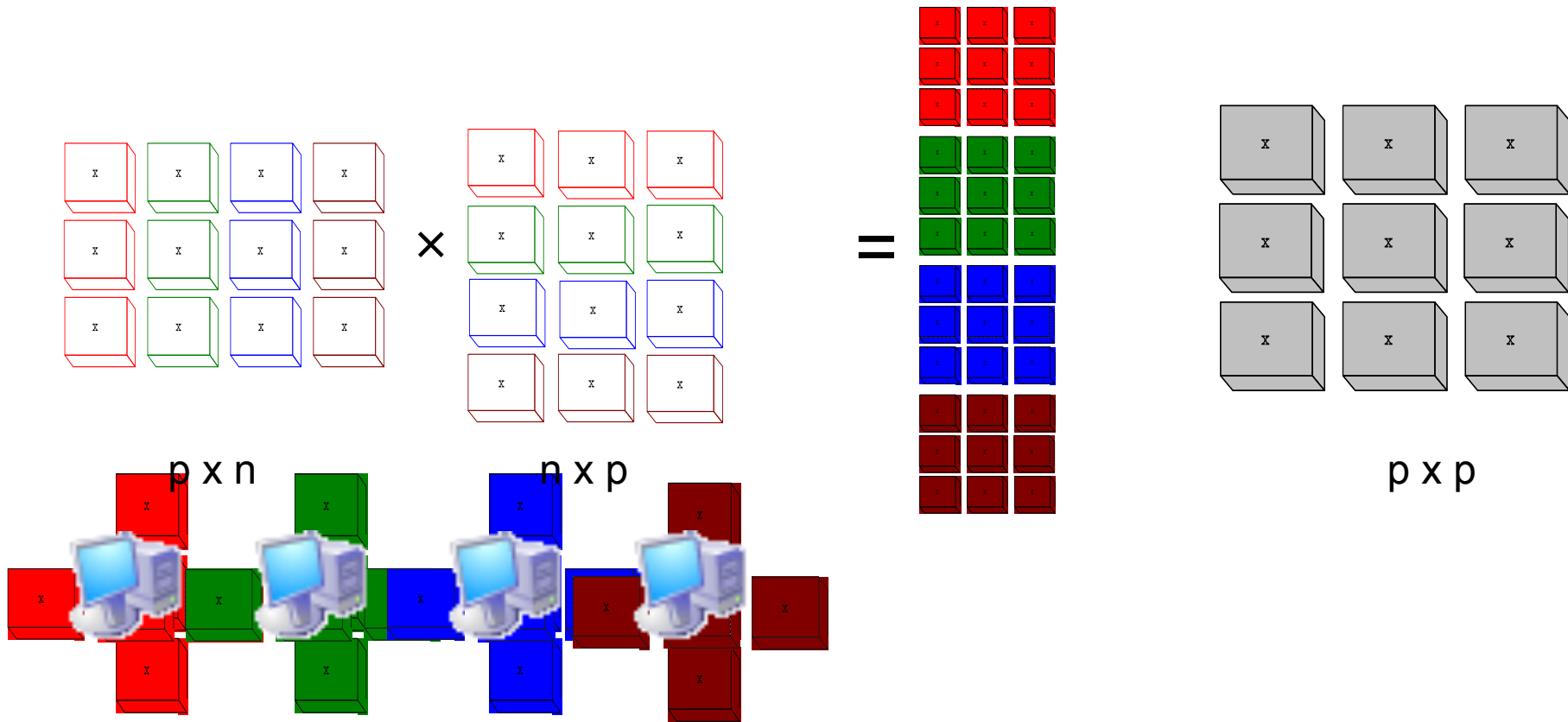




# PSVM



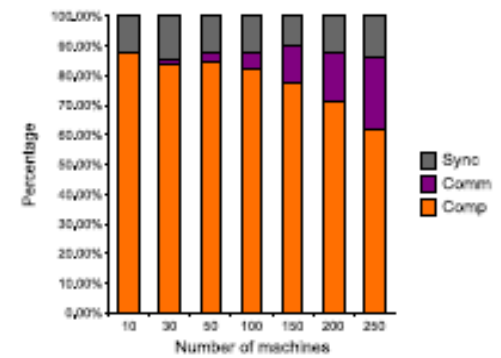
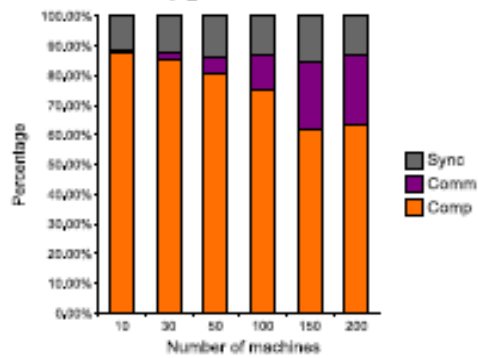
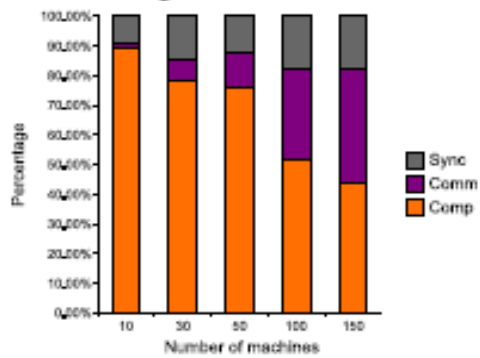
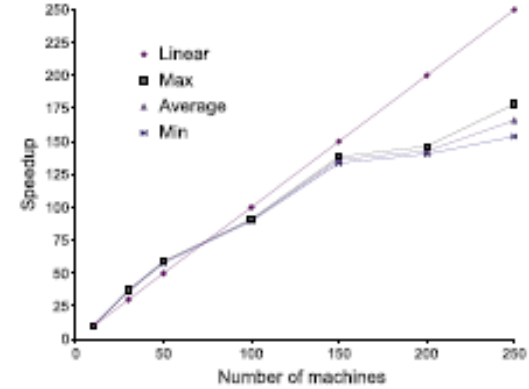
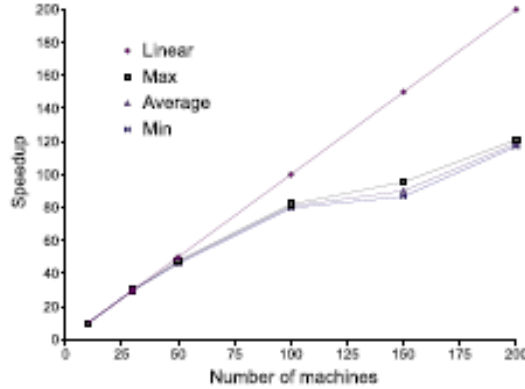
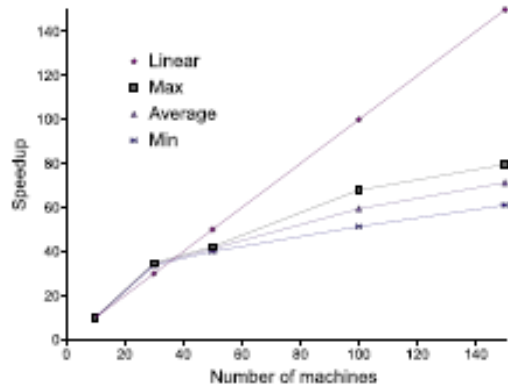
# Matrix Product



# PSVM [E. Chang, et al, NIPS 07]

- Column-based ICF
  - Slower than row-based on single machine
  - Parallelizable on multiple machines
- Changing IPM computation order to achieve parallelization

# Overheads



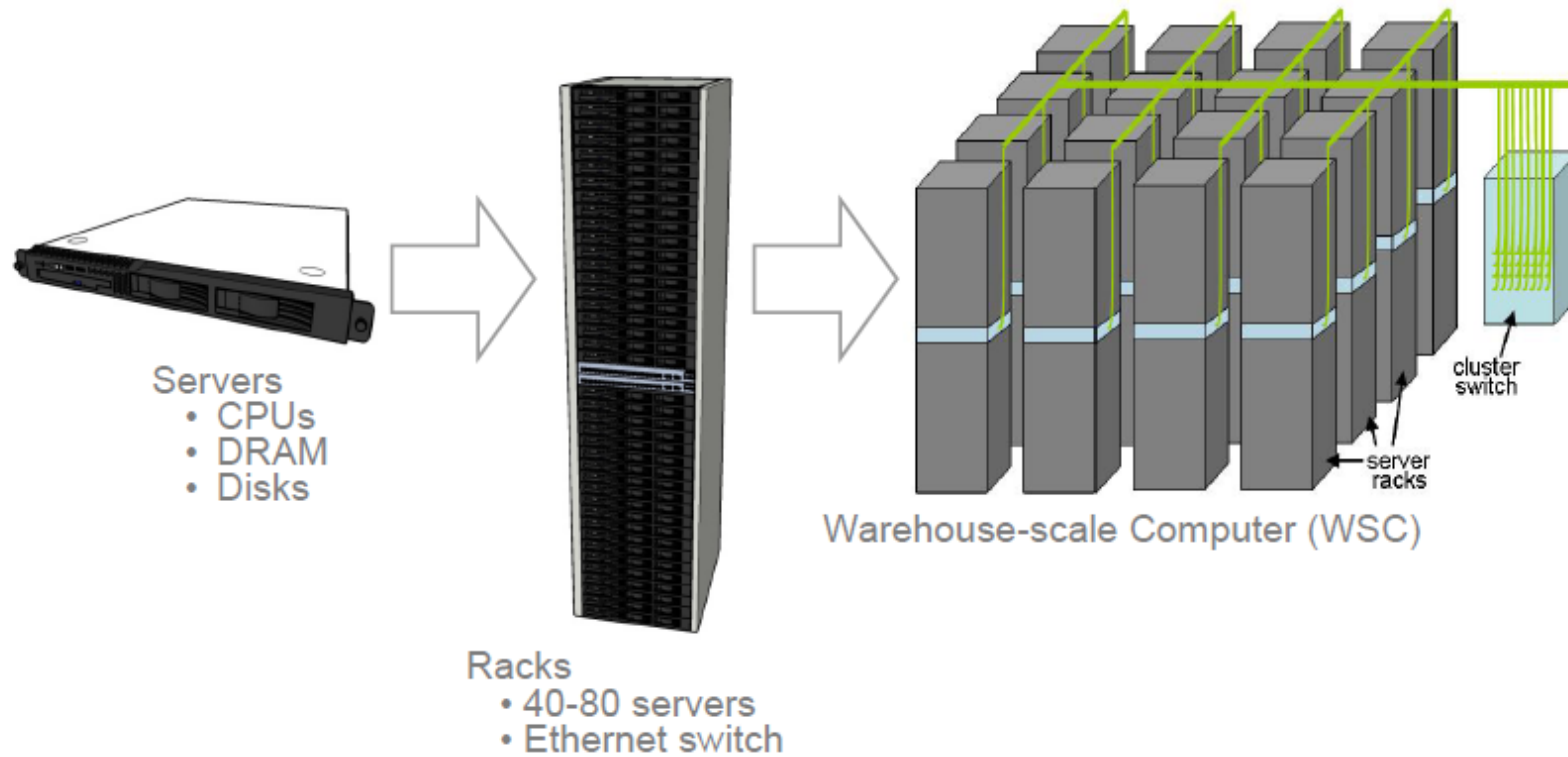
# Speedup

Machines	Image (200k)		CoverType (500k)		RCV (800k)	
	Time (s)	Speedup	Time (s)	Speedup	Time (s)	Speedup
10	1,958 (9)	10*	16,818 (442)	10*	45,135 (1373)	10*
30	572 (8)	34.2	5,591 (10)	30.1	12,289 (98)	36.7
50	473 (14)	41.4	3,598 (60)	46.8	7,695 (92)	58.7
100	330 (47)	59.4	2,082 (29)	80.8	4,992 (34)	90.4
150	274 (40)	71.4	1,865 (93)	90.2	3,313 (59)	136.3
200	294 (41)	66.7	1,416 (24)	118.7	3,163 (69)	142.7
250	397 (78)	49.4	1,405 (115)	119.7	2,719 (203)	166.0
500	814 (123)	24.1	1,655 (34)	101.6	2,671 (193)	169.0
LIBSVM	4,334 NA	NA	28,149 NA	NA	184,199 NA	NA

# Scalability

- Computation
  - Parallelization
  - Approximation
- File Systems
  - Latency
  - Recovery
- Power Management

# Sample Platforms

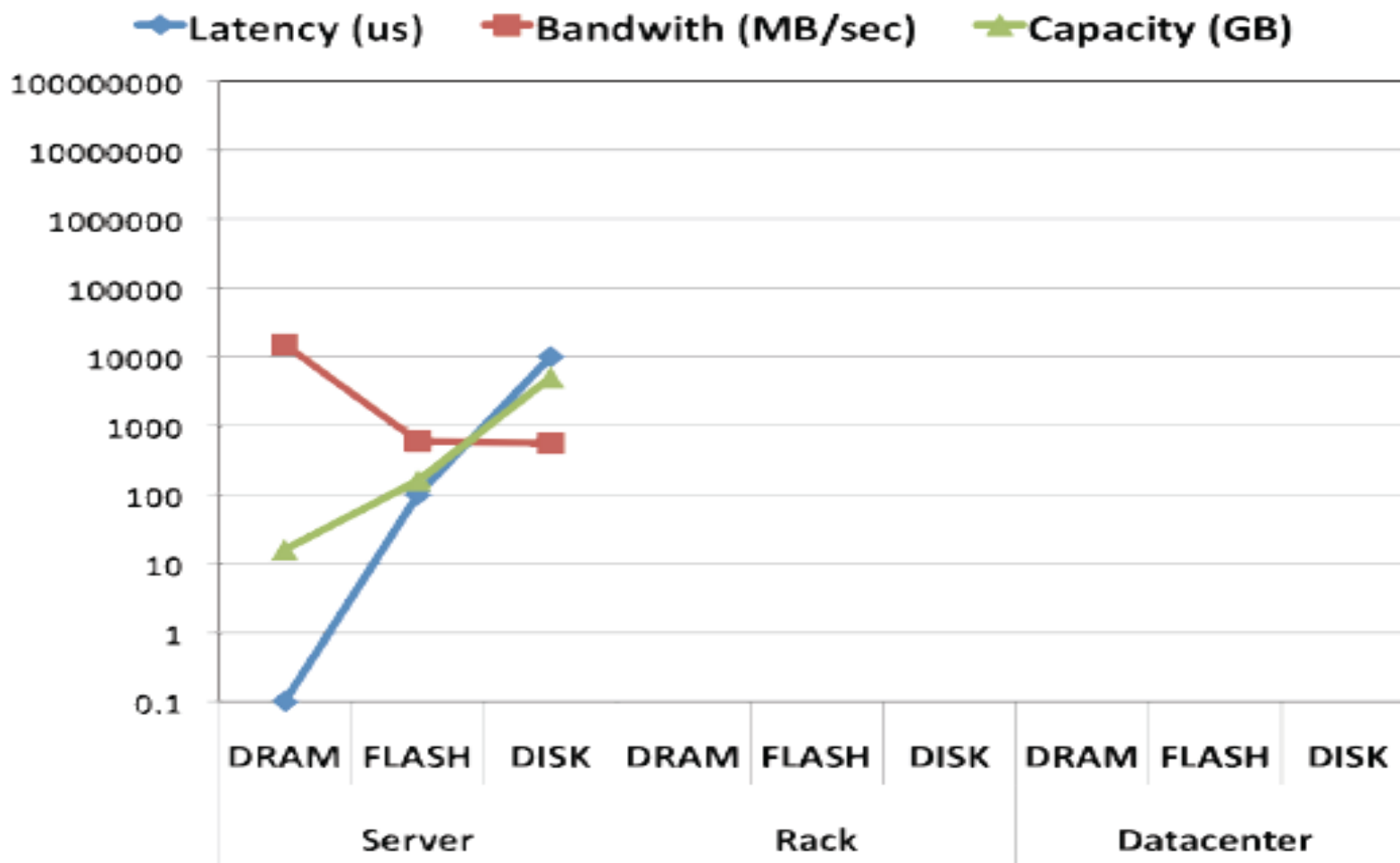


# Sample Hierarchy

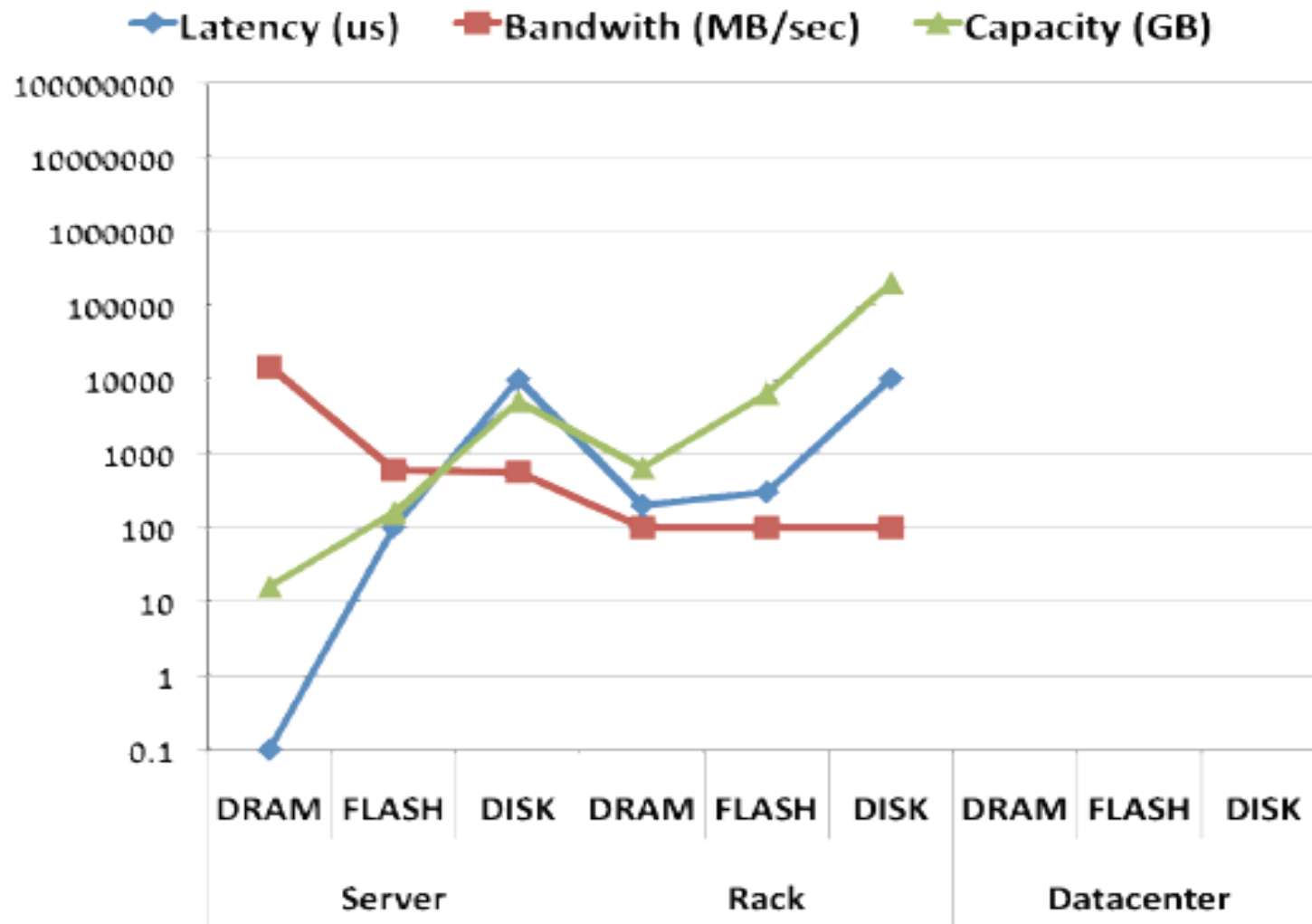
- Server
  - 16GB DRAM; 160MB Flash; 5 x 1TB disk
- Rack
  - 40 servers
  - 48 port Gigabit Ethernet switch
- Warehouse
  - 10,000 servers (250 racks)
  - 2K port Gigabit Ethernet switch



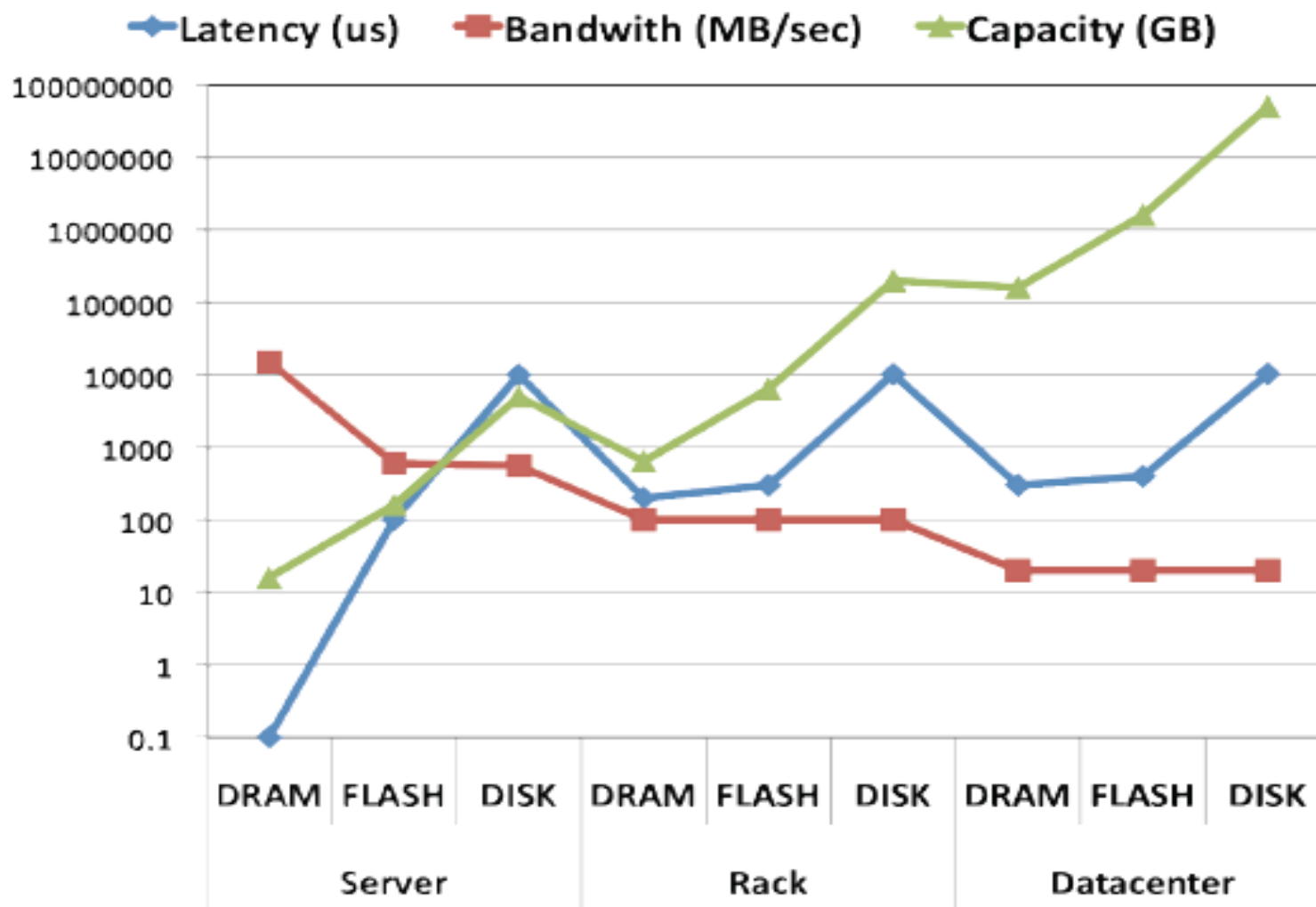
# Storage --- One Server



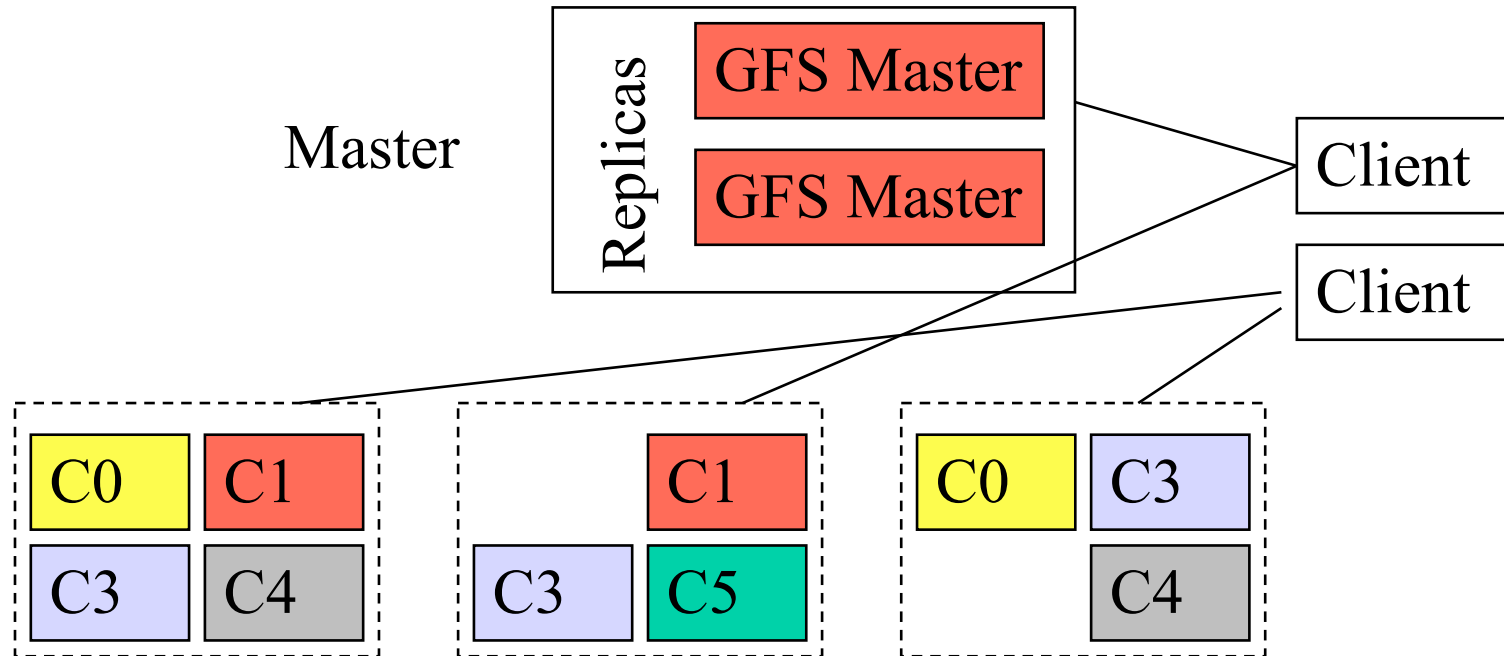
# Storage --- One Rack



# Storage --- One Center



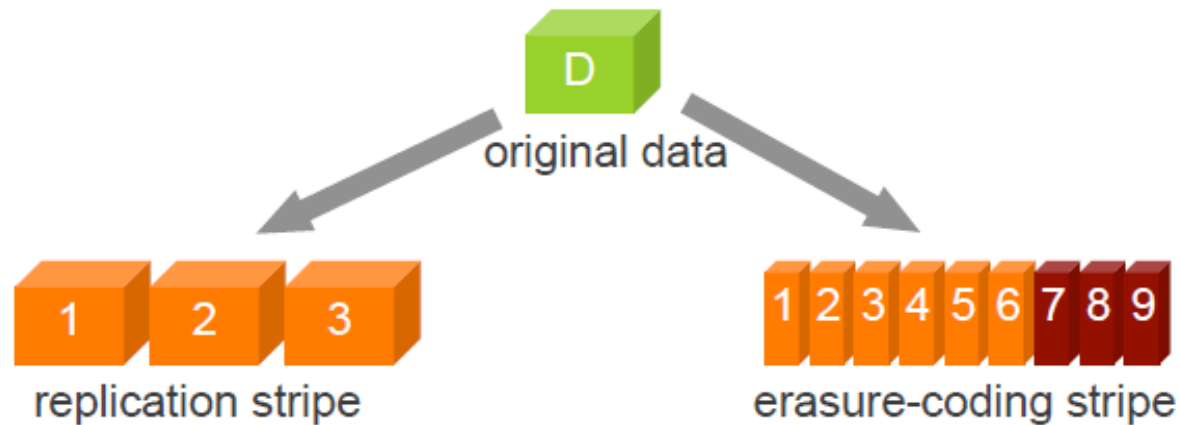
# Google File System (GFS)



- Master manages metadata
- Data transfers happen directly between clients/chunkservers
- Files broken into chunks (typically 64 MB)
- Chunks triplicated across three machines for safety
- See SOSP<sup>03</sup> paper at <http://labs.google.com/papers/gfs.html>

# WSC data availability: cluster file systems

---



- Data blocks of each stripe are placed on different fault domains
  - different disks, servers, racks
  - Data blocks are distributed across the whole WSC
    - read operations are easily load-balanced
    - recovery is highly efficient
- What affects data availability as seen by a client of a cluster file system?

# Win in Scale

- Google Translate
- Voice
- Trend Prediction
  - An example benefits society

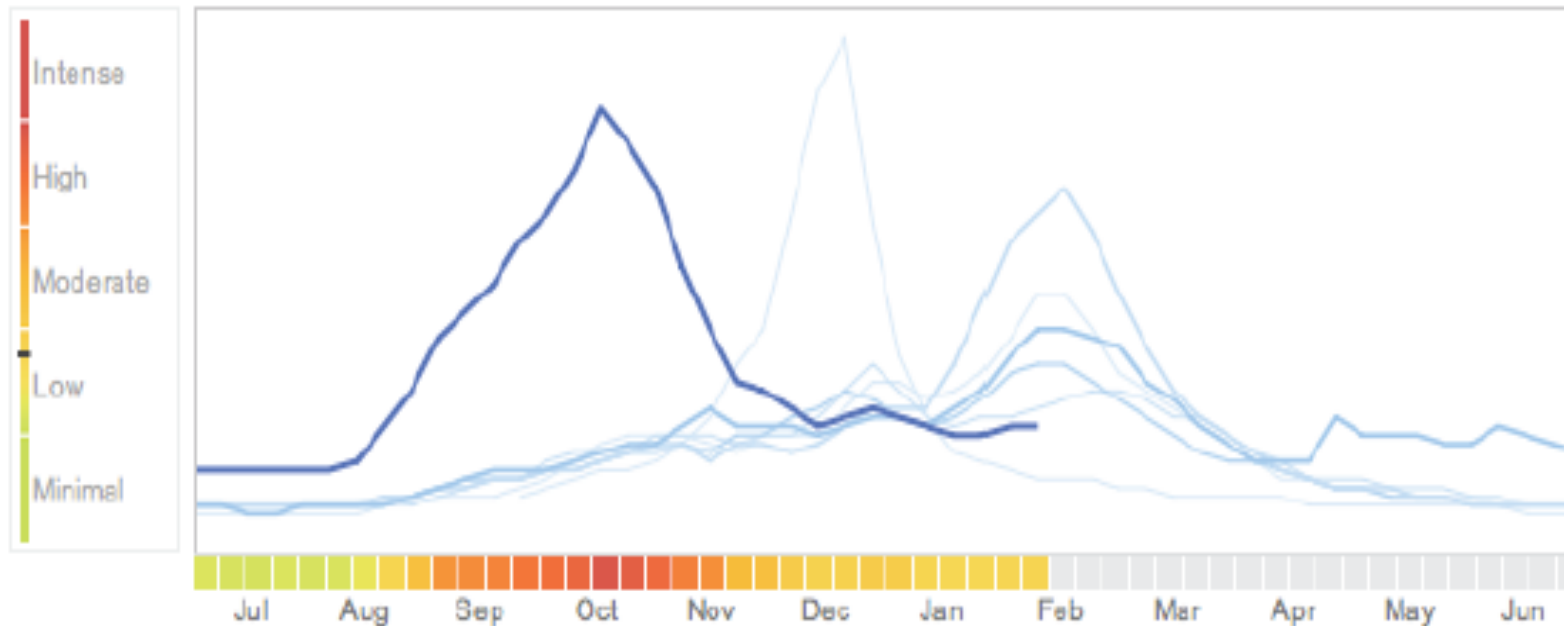
# H1N1 United Nation Report

## Explore flu trends - United States

We've found that certain search terms are good indicators of flu activity. Google Flu Trends uses aggregated Google search data to estimate flu activity. [Learn more »](#)

### National

● 2009-2010 ● [Past years ▼](#)



# Concluding Remarks

- Search + Social
- Increasing quantity and complexity of data demands scalable solutions
- Have parallelized key subroutines for mining massive data sets
  - Spectral Clustering [ECML 08]
  - Frequent Itemset Mining [ACM RS 08]
  - PLSA [KDD 08]
  - LDA [WWW 09, AAIM 09]
  - UserRank [Google TR 09]
  - Support Vector Machines [NIPS 07]
- Launched Google Q&A (Confucius) in 60+ countries
- Relevant papers
  - <http://infolab.stanford.edu/~echang/>
- Open Source PSVM, PLDA
  - <http://code.google.com/p/psvm/>
  - <http://code.google.com/p/plda/>



# Models of Innovation

- Ivory tower
  - Only consider theory but not application
- Build it and they will come
  - Scientists drives product development
- “Research for sale”
  - Research funded by:
    - product groups or customers
- Research & development as equals
  - Research “sells” innovation;
  - Product “requests” innovation
- Google-style innovation

