

Two Random Walks that Surprise

Ashish Goel
Stanford University

PageRank




Pankaj Gupta
@pankaj

TWEETS **5,520** FOLLOWING **949** FOLLOWERS **9,329**

Compose new Tweet...

Who to follow · Refresh · View all

-  **Safeway** @Safeway
Followed by Mazen Rawashd...
Follow Promoted
-  **Discovery News** @DNews
Follow
-  **Micheal** @micheal
Followed by Vinod Kone an...
Follow


Popular accounts · Find friends

Trends · Change


- #StPatricksDay
- L'Wren Scott
- #NCAA
- Venezuela
- Jim Irsay
- DRC
- #mcm
- #luckoftheirish
- #startups
- Modi

Tweets

3 new Tweets




Aapo Kyrola @kyrpov · 1m
GraphChi used for computational biology: homes.di.unimi.it/valentini/pape...
Results: As fast as in-memory computation, way faster than Neo4j.
Expand Reply Retweet Favorite Pocket More



Rohit @rohit_x_ · 1m
This is big. theguardian.com/science/2014/m... Primordial Gravitational waves 'seen' in polarization of early light. #Einstein
View summary Reply Retweet Favorite Pocket More



D-Lab @ MIT @dlab_mit · 2m
March 26: Follow-up @harvest_fuel webinar to March 5th "Charcoal Briquette Enterprise Development" @dlab_mit @harvest_fuel...
Expand Reply Retweet Favorite More



Mark McBride @mccv · 7m
today I'm responding to all iMessages with "congratulations!".
Expand Reply Retweet Favorite More

1 more reply



Mark McBride @mccv · 4m
@matasar you know that's like me clowning you for the links in DMs debacle
Expand Reply Retweet Favorite More



Ben Matasar @matasar · 3m
@mccv So?
Expand Reply Retweet Favorite More



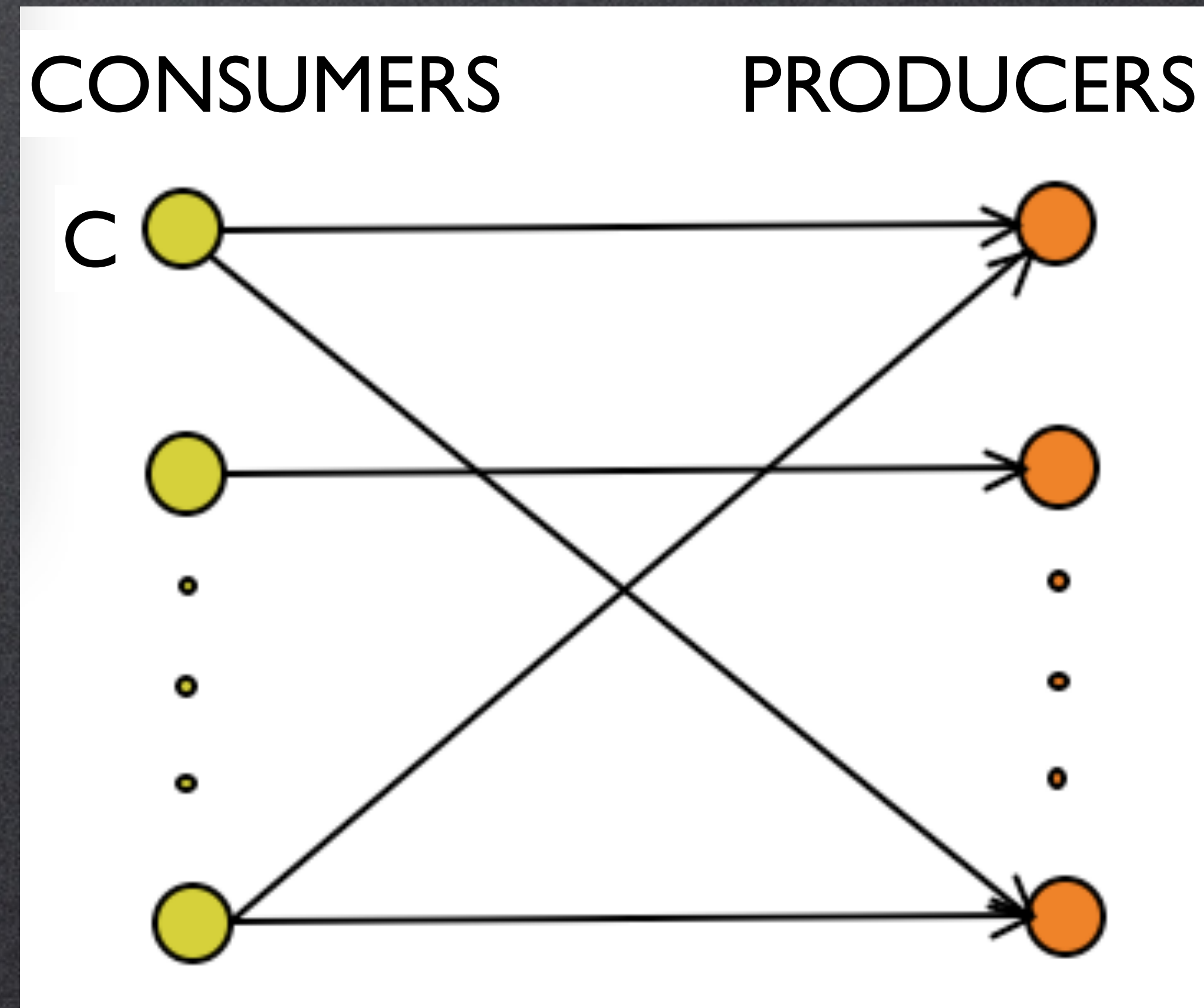
Eater SF @eatersf · 6m
Check out the scene at @offthegridsf's Fort Mason opening night, see what's new this year. eater.cc/1iwevsR pic.twitter.com/cjSXQMdOdn



Collaborative Filtering

To get recommendations for C, compute similarity scores for all consumers, and relevance scores for all producers, with respect to C

1. Start with $\text{sim}(C) = 1$
2. Propagate similarity scores along graph edges to compute relevance scores, and vice-versa



Many propagation methods; Often, a linear system of equations

Collaborative Filter: Love or Money

How should we do this propagation? Two extremes:

LOVE: All the similarity score of a consumer X gets transferred to **each** producer that X follows, and the same in the reverse direction

→ Analogous to Singular Value Decompositions in the dense graph limit (HITS)

MONEY: If X follows d producers, then a fraction $1/d$ of the similarity score of X gets transferred to each producer that X follows (SALSA)

Personalized PageRank

Given a **consumer C**, perform a random walk on the Follow graph. If the walk is at node **v**, then the walk:

- Jumps back to node **C** with probability α
- Follows a random edge out of **v** with probability $1 - \alpha$

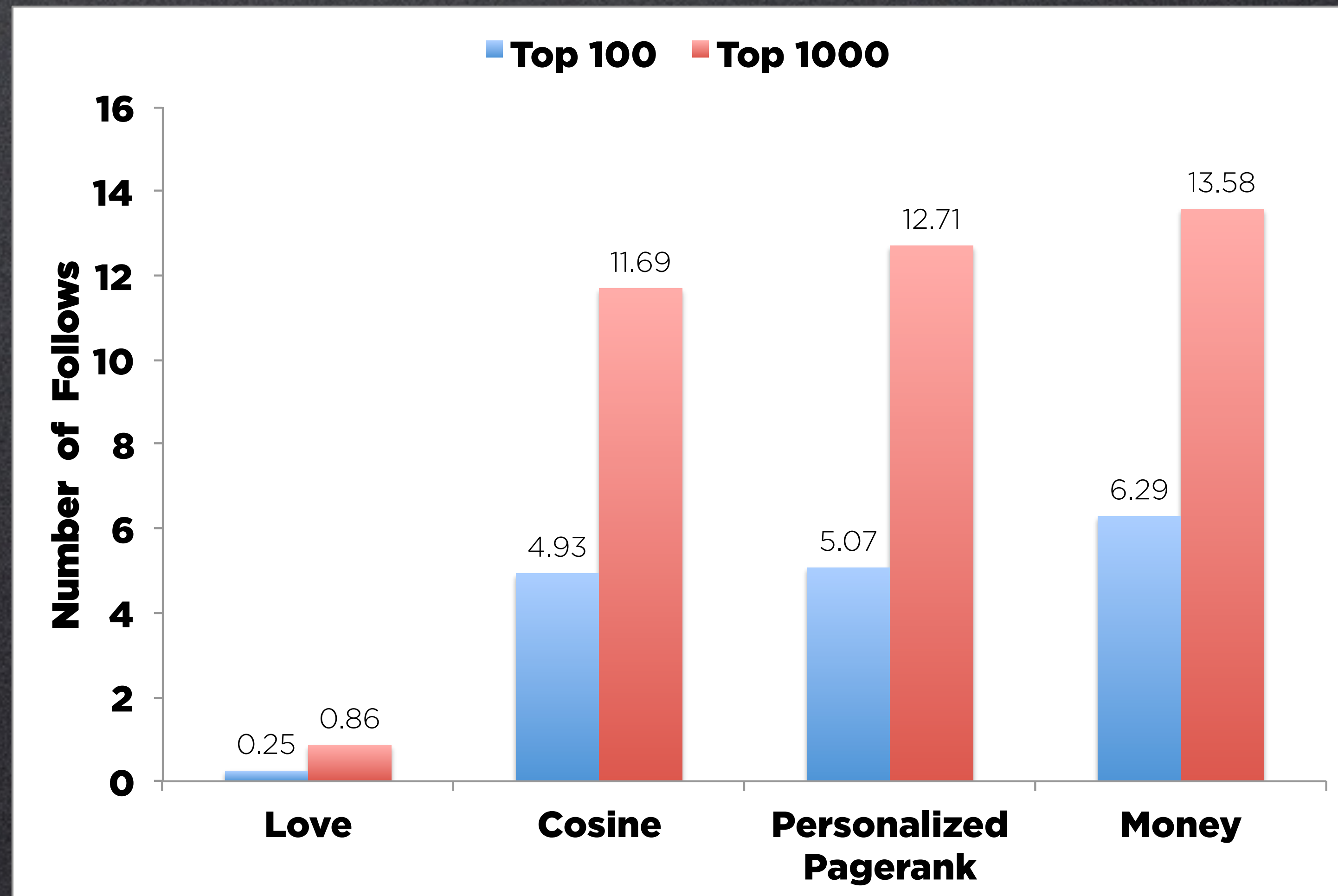
The Personalized PageRank of node **Y** is the weight of **Y** in the stationary distribution of this random walk

SALSA/Money is just Personalized PageRank run on the undirected consumer–producer graph

A Dark Test

Run various algorithms to predict follows, but don't display the results. Instead, just observe how many of the top **predictions** get followed organically

[Bahmani, Chowdhury, Goel; 2010]





Strategic Impact

Creates billions of new follows every year

- More than 1/8 of new follows are directly via the Who-to-Follow module
- More than 15% of active users (> 36 Million users) make at least one follow every month via this module

Promoted Tweets and Promoted Accounts

The screenshot displays the Twitter interface with a navigation bar at the top containing icons for Home, Notifications, Search, Profile, and the Twitter logo. The main content is divided into two columns.

Left Column: Who to follow

- Who to follow · Refresh · View all**
- CKM Advisors** @CKMAdvi...
 Follow Promoted
- girish sastry** @girishsastry
 Followed by Utkarsh Srivast...
 Follow
- Shiv Ramamurthi** @mogro...
 Followed by Stanford Alumn...
 Follow
- Popular accounts · Find friends

Right Column: Tweets

1 new Tweet

- Aneesh Sharma** @aneeshs · 4m
 Feeling lucky to be at #analytics2014 with @ashishgoel @johnsirois @pankaj @sgurumur for our #edelmanaward presentation. Go #teamtwitter!
 Expand
 Reply Retweet Favorite More
- John Sirois** @johnsirois · 5m
 Hanging out with @ashishgoel @sgurumur @pankaj @aneeshs #analytics2014. Special thanks to our #edelmanaward coaches John Birge & Carrie Beam
 Expand
 Reply Retweet Favorite More
- Followed by Peter Fenton.
 NewRelic @newrelic · Mar 11
 4 Essential Tips from the Coding CEO. How New Relic CEO Lew Cirne still builds product: blog.newrelic.com/2014/03/11/sxs...
 Promoted by NewRelic
 Expand
 Reply Retweet Favorite More

Promoted Tweets and Promoted Accounts

The screenshot displays the Twitter interface with a navigation bar at the top containing icons for home, notifications, search, profile, and the Twitter logo. On the left side, there is a 'Who to follow' section with three accounts: CKM Advisors (promoted), girish sastry, and Shiv Ramamurthi. On the right side, the 'Tweets' section shows a list of tweets. The third tweet, from NewRelic, is circled in purple. This tweet is promoted and includes a link to a blog post. The interface also shows navigation icons at the bottom.

Who to follow · Refresh · View all

- CKM Advisors** @CKMAdvi... **Follow** Promoted
- girish sastry** @girishsastry **Follow** Followed by Utkarsh Srivast...
- Shiv Ramamurthi** @mogro... **Follow** Followed by Stanford Alumn...

Popular accounts · Find friends

Tweets

1 new Tweet

- Aneesh Sharma** @aneeshs · 4m
Feeling lucky to be at #analytics2014 with @ashishgoel @johnsirois @pankaj @sgurumur for our #edelmanaward presentation. Go #teamtwitter!
Expand Reply Retweet Favorite More
- John Sirois** @johnsirois · 5m
Hanging out with @ashishgoel @sgurumur @pankaj @aneeshs #analytics2014. Special thanks to our #edelmanaward coaches John Birge & Carrie Beam
Expand Reply Retweet Favorite More
- Followed by Peter Fenton.
NewRelic @newrelic · Mar 11
4 Essential Tips from the Coding CEO. How New Relic CEO Lew Cirne still builds product: blog.newrelic.com/2014/03/11/sxs...
Promoted by NewRelic
Expand Reply Retweet Favorite More

Promoted Tweets and Promoted Accounts

The image shows a screenshot of the Twitter mobile app interface. On the left, the 'Who to follow' section is visible, featuring three account suggestions. The first, 'CKM Advisors @CKMAdvi...', is highlighted with a purple oval and includes a 'Promoted' badge. Below it are 'gIrish sastry @gIrishsastry' and 'Shiv Ramamurthi @mogro...'. At the bottom of this section are links for 'Popular accounts' and 'Find friends'. On the right, the 'Tweets' section shows a '1 new Tweet' notification. The first tweet is from 'Aneesh Sharma @aneeshs' about the #analytics2014 event. The second tweet is from 'John Sirois @johnsirois' about the same event. A third tweet from 'NewRelic @newrelic' is also highlighted with a purple oval; it includes a 'Followed by Peter Fenton.' notification and a 'Promoted by NewRelic' badge. The tweet text discusses '4 Essential Tips from the Coding CEO' and includes a link to a blog post.

Who to follow · Refresh · View all

CKM Advisors @CKMAdvi... Follow Promoted

gIrish sastry @gIrishsastry Follow
Followed by Utkarsh Srivast...

Shiv Ramamurthi @mogro... Follow
Followed by Stanford Alumn...

Popular accounts · Find friends

Tweets

1 new Tweet

Aneesh Sharma @aneeshs · 4m
Feeling lucky to be at #analytics2014 with @ashishgoel @johnsirois @pankaj @sgurumur for our #edelmanaward presentation. Go #teamtwitter!
Expand Reply Retweet Favorite More

John Sirois @johnsirois · 5m
Hanging out with @ashishgoel @sgurumur @pankaj @aneeshs #analytics2014. Special thanks to our #edelmanaward coaches John Birge & Carrie Beam
Expand Reply Retweet Favorite More

Followed by Peter Fenton.

NewRelic @newrelic · Mar 11
4 Essential Tips from the Coding CEO. How New Relic CEO Lew Cirne still builds product: blog.newrelic.com/2014/03/11/sxs...
Promoted by NewRelic
Expand Reply Retweet Favorite More

Impact on Revenue

“The Who-To-Follow system was crucial, in a fundamental way, for the Promoted Accounts product, and the Promoted Tweets product also initially used the Who-To-Follow system’s targeting”

– Alex Roetter (VP of Engineering, Revenue)

Scientific Questions

1. Fast Incremental PageRank
2. Fast Personalized PageRank

Incremental PageRank

Updates to social graph are made in real-time

- As opposed to a batched crawl process for web search
- Real-time updates to PageRank are important to capture trending events

Goal: Design an algorithm to update PageRank incrementally (i.e. upon an edge arrival)

- t -th edge arrival: Let (u_t, v_t) denote the arriving edge, $d_t(v)$ denote the out-degree of node v , and $\pi_t(v)$ its PageRank

Incremental PageRank via Monte Carlo

Start with $R = O(\log N)$ random walks from every node

At time t , for every random walk through node u_t , re-route it to use the new edge (u_t, v_t) with probability $1/d_t(u_t)$

→ Time/number of network-calls for each re-routing: $O(1/\alpha)$

Claim: This faithfully maintains R random walks after arbitrary edge arrivals

Need the graph and the stored random walks in fast distributed memory

Incremental PageRank Time

Assume that the edges of the graph are chosen by an adversary, but then presented in random order

Theorem: # of re-routings per arrival goes to 0

→ t -th arrival: # of reroutes = $O(N R / (\alpha t))$

→ Total time over M arrivals = $O((N R \log N) / \alpha^2)$

→ Comparable to doing power iteration/Monte Carlo just once!

[Bahmani, Goel, Chowdhury, VLDB 2010]

Incremental PageRank Time

Assume that the edges of the graph are chosen by an adversary, but then presented in random order

Theorem: # of **Monte Carlo** per arrival goes to 0

→ t-th arrival: # of **Monte Carlo** takes time $NR/(\alpha t)$

→ Total time over M arrivals = $O((NR \log N)/\alpha^2)$

→ Comparable to doing power iteration/**Monte Carlo** just once!

[Bahmani, Goel, Chowdhury, VLDB 2010]

Incremental PageRank Time

Assume that the edges of the graph are chosen by an adversary, but then presented in random order

Theorem: # of **Monte Carlo** per arrival goes to 0

→ t-th arrival: # of **Monte Carlo** takes time $NR/(\alpha t)$

→ Total time over M arrivals = $O((NR \log N)/\alpha^2)$

→ Comparable to doing power iteration/**Monte Carlo** just once!

**Only an extra
log N/α**

[Bahmani, Goel, Chowdhury, VLDB 2010]

Incremental PageRank Time

Assume that the edges of the graph are chosen by an adversary, but then presented in random order

Theorem: # of **Power iteration takes time $M R/\alpha$** per arrival goes to 0

→ t-th arrival: # of **Power iteration takes time $M R/\alpha$** $N R/(\alpha t)$

→ Total time over M arrivals = $O((N R \log N)/\alpha^2)$

→ Comparable to doing **power iteration**/Monte Carlo just once!

[Bahmani, Goel, Chowdhury, VLDB 2010]

Incremental PageRank Time

Assume that the edges of the graph are chosen by an adversary, but then presented in random order

Theorem: # of **power iteration** per arrival goes to 0

→ t-th arrival: # of **power iteration** $N R / (\alpha t)$

→ Total time over M arrivals = $O((N R \log N) / \alpha^2)$

→ Comparable to doing **power iteration**/Monte Carlo just once!

Power iteration takes time $M R / \alpha$

$N \log N / \alpha$ vs M

[Bahmani, Goel, Chowdhury, VLDB 2010]

Personalized PageRank

Network-based Personalized Search is not yet mature

Missing technical piece: Efficient algorithms for Personalized PageRank Queries

→ Given source s and target t , estimate the Personalized PageRank of t for s with high accuracy, if it is greater than δ

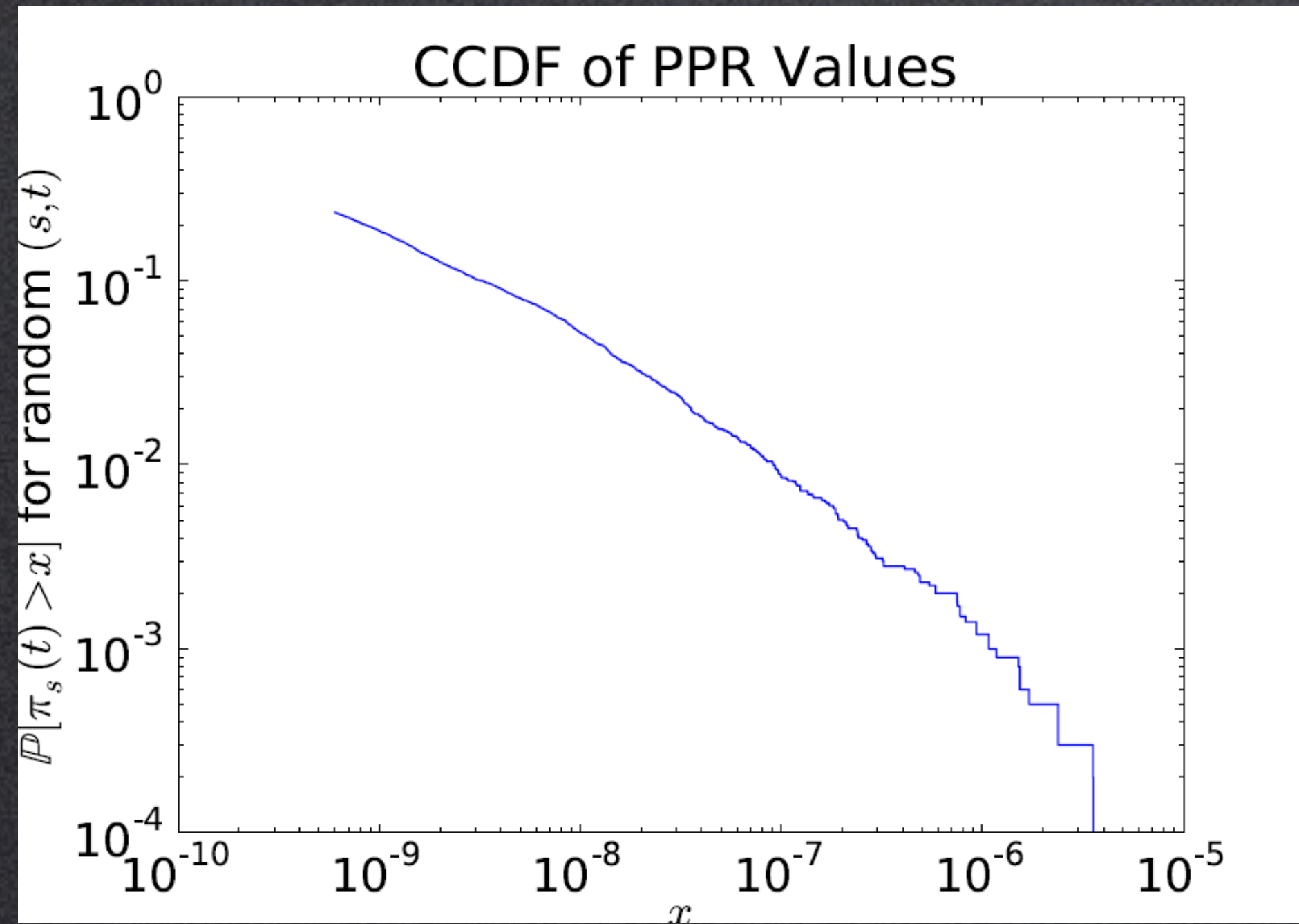
Personalized PageRank

Given a **consumer C**, perform a random walk on the Follow graph. If the walk is at node v , then the walk:

- Jumps back to node C with probability α
- Follows a random edge out of v with probability $1 - \alpha$

The Personalized PageRank of node Y is the weight of Y in the stationary distribution of this random walk

Existing Methods for PPR Queries



Monte Carlo uses time $> 1/\delta$
“Local Update” uses time d/δ

[$d = M/N$ is the average degree]

On Twitter-2010, if $\delta = \frac{4}{n} \approx 10^{-7}$, then

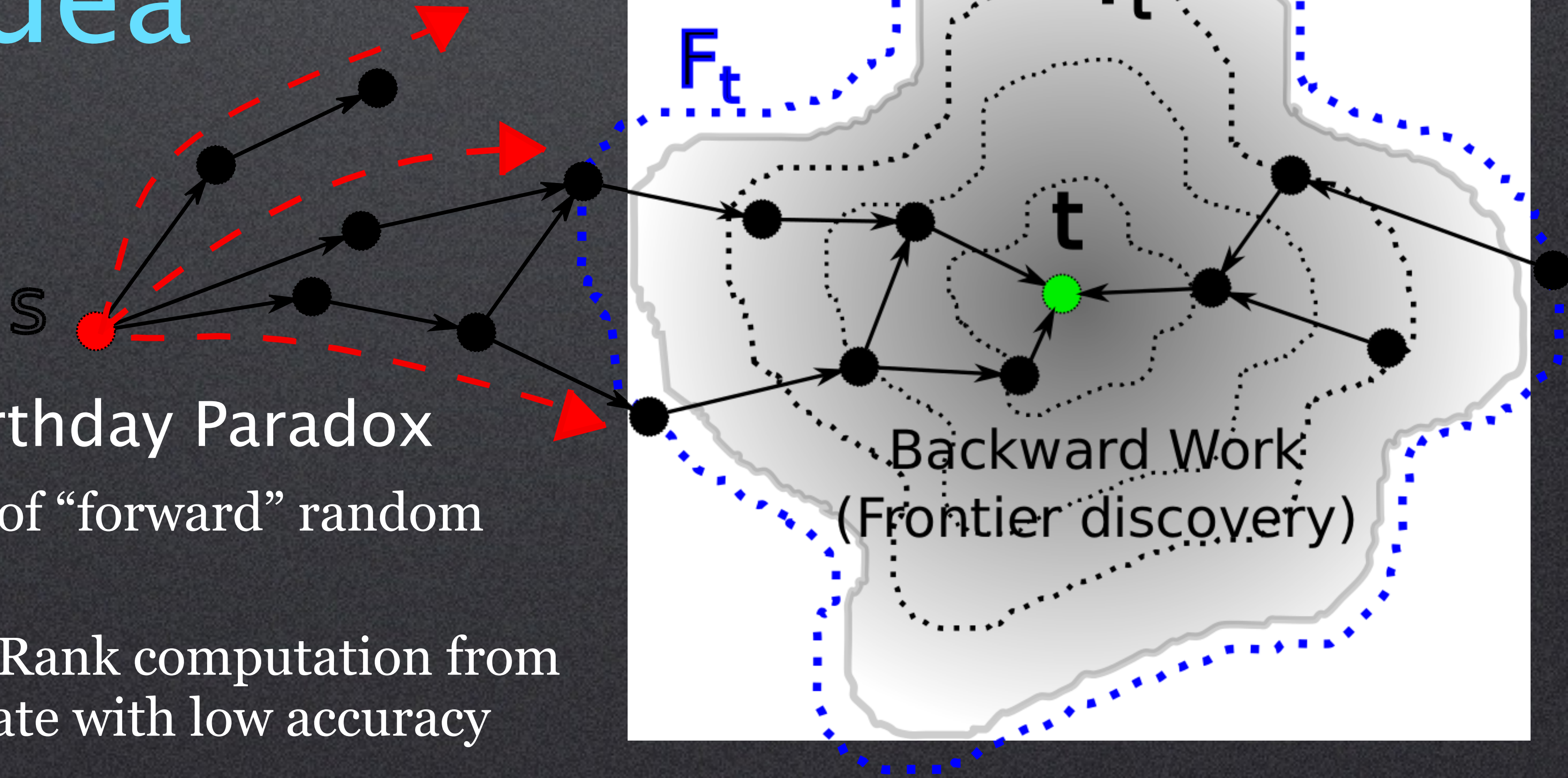
$$\Pr[\pi(s, t) > \delta] = 1\%$$

FAST PPR

We can answer PPR queries in either

- Average time $\tilde{O}(\sqrt{d/\delta})$
- Worst case time $\tilde{O}(\sqrt{d/\delta})$ with $\tilde{O}(\sqrt{d/\delta})$ storage and pre-processing time per node
- Typical values: $\delta \sim 10^{-8}$, $d \sim 100$; results in a > 100 -fold decrease

Basic Idea



Intuition: The Birthday Paradox

- Do small number of “forward” random walks from s
- Do “reverse” PageRank computation from t using Local Update with low accuracy
- Use number of collisions as an estimator
- Need to “catch” a collision just before it happens

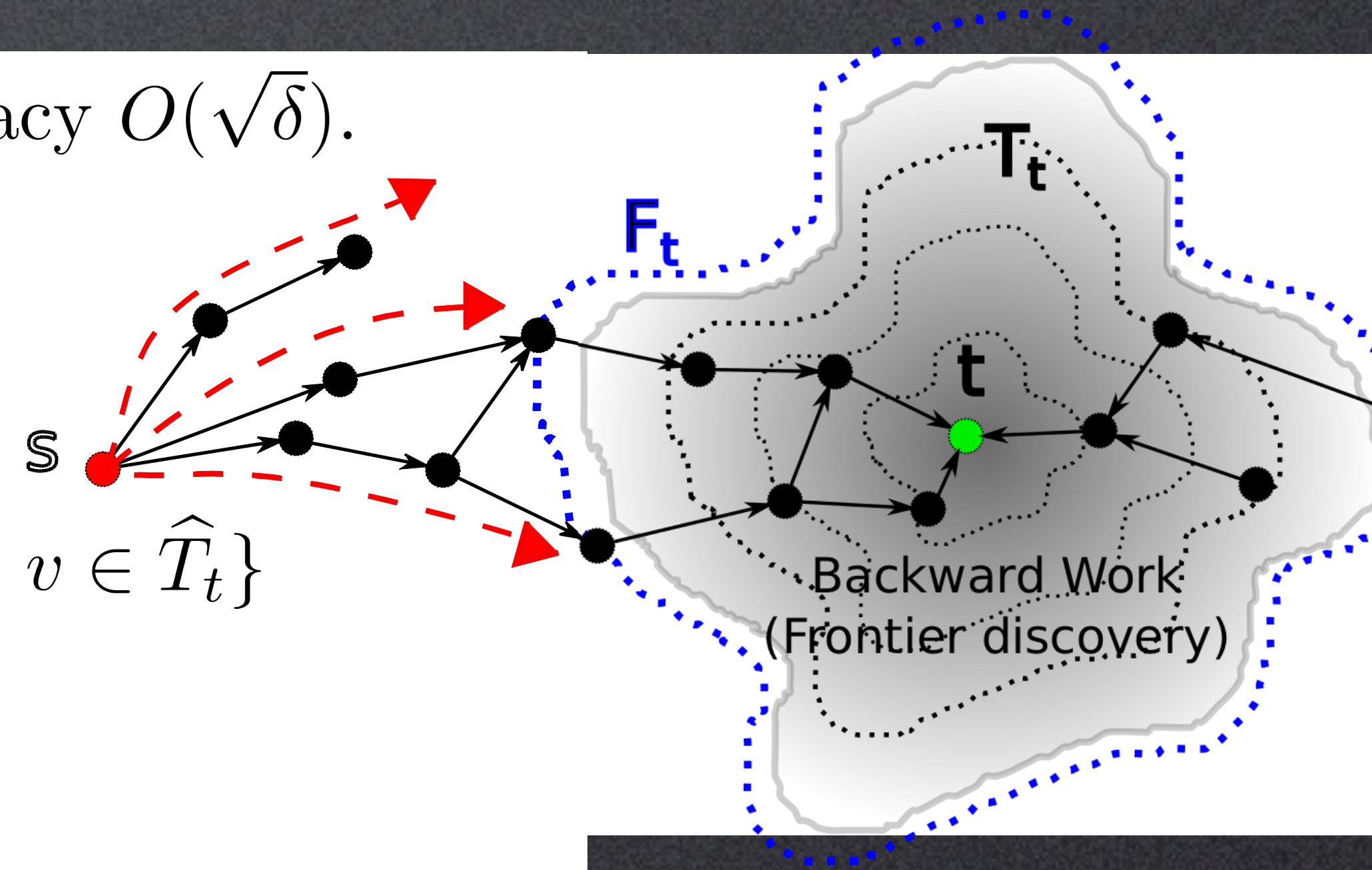
Simple Version of FAST PPR

1. Use Local Update to compute estimates $\hat{\pi}(v, t)$ to accuracy $O(\sqrt{\delta})$.

2. Define

$$\text{Target Set } \hat{T}_t = \{v \in V : \hat{\pi}(v, t) > \sqrt{\delta}\}$$

$$\text{Frontier } \hat{F}_t = \{u \in V \setminus \hat{T}_t : (u, v) \in E \text{ for some } v \in \hat{T}_t\}$$



3. Take $O\left(\frac{\log(n)}{\sqrt{\delta}}\right)$ Random Walks $\{W_i\}$, terminating each early if it hits \hat{F}_t .

Define

$$X_i = \begin{cases} \hat{\pi}(u, t), & W_i \text{ hits } u \in \hat{F}_t \\ 0, & W_i \text{ does not hit } \hat{F}_t \end{cases}$$

4. Return empirical mean $\{X_i\}$.

Running Time for Simple Version

For a uniformly random target node t , the average per-query running time is

$$O\left(\frac{1}{\sqrt{\delta}} (\bar{d} + \log(n))\right).$$

Running Time for Simple Version

For a uniformly random target node t , the average per-query running time is

$$O\left(\frac{1}{\sqrt{\delta}} (\bar{d} + \log(n))\right).$$

Reverse work
(Local Update)



Running Time for Simple Version

For a uniformly random target node t , the average per-query running time is

$$O\left(\frac{1}{\sqrt{\delta}} (\bar{d} + \log(n))\right).$$

Reverse work
(Local Update)



Forward work
(Monte Carlo)



Running Time for Simple Version

For a uniformly random target node t , the average per-query running time is

$$O\left(\frac{1}{\sqrt{\delta}}(\bar{d} + \log(n))\right).$$

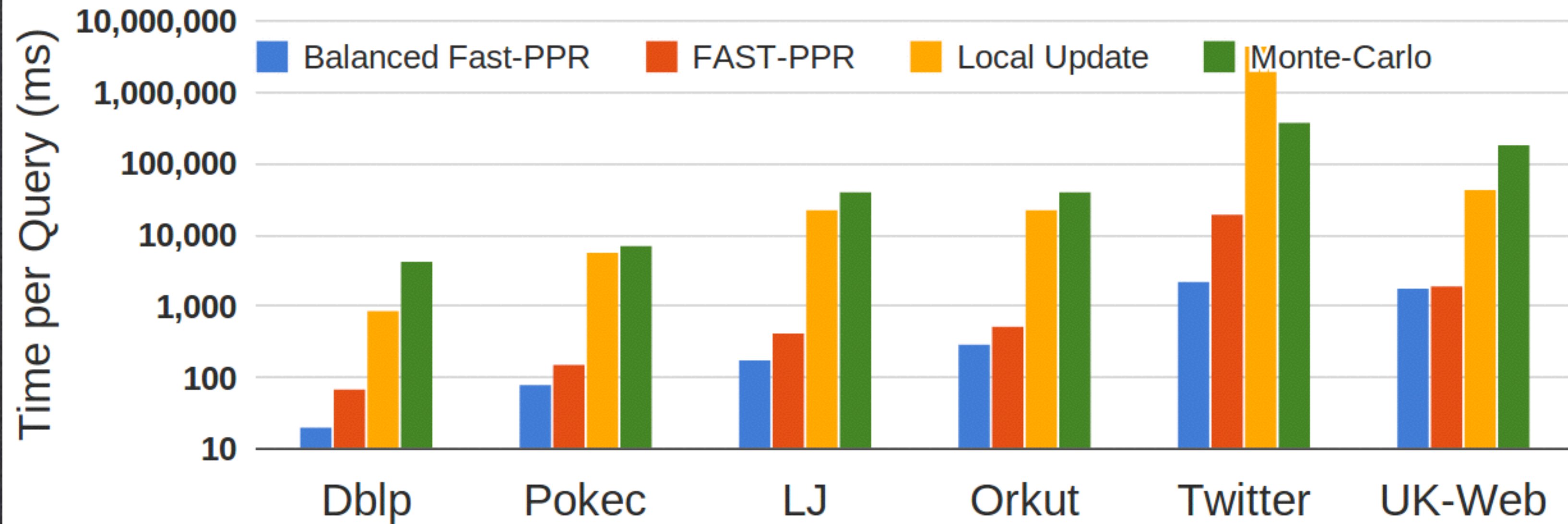
Reverse work
(Local Update)

Forward work
(Monte Carlo)

We get final running time of $\tilde{O}(\sqrt{d/\delta})$ by using different accuracies in forward and reverse computation

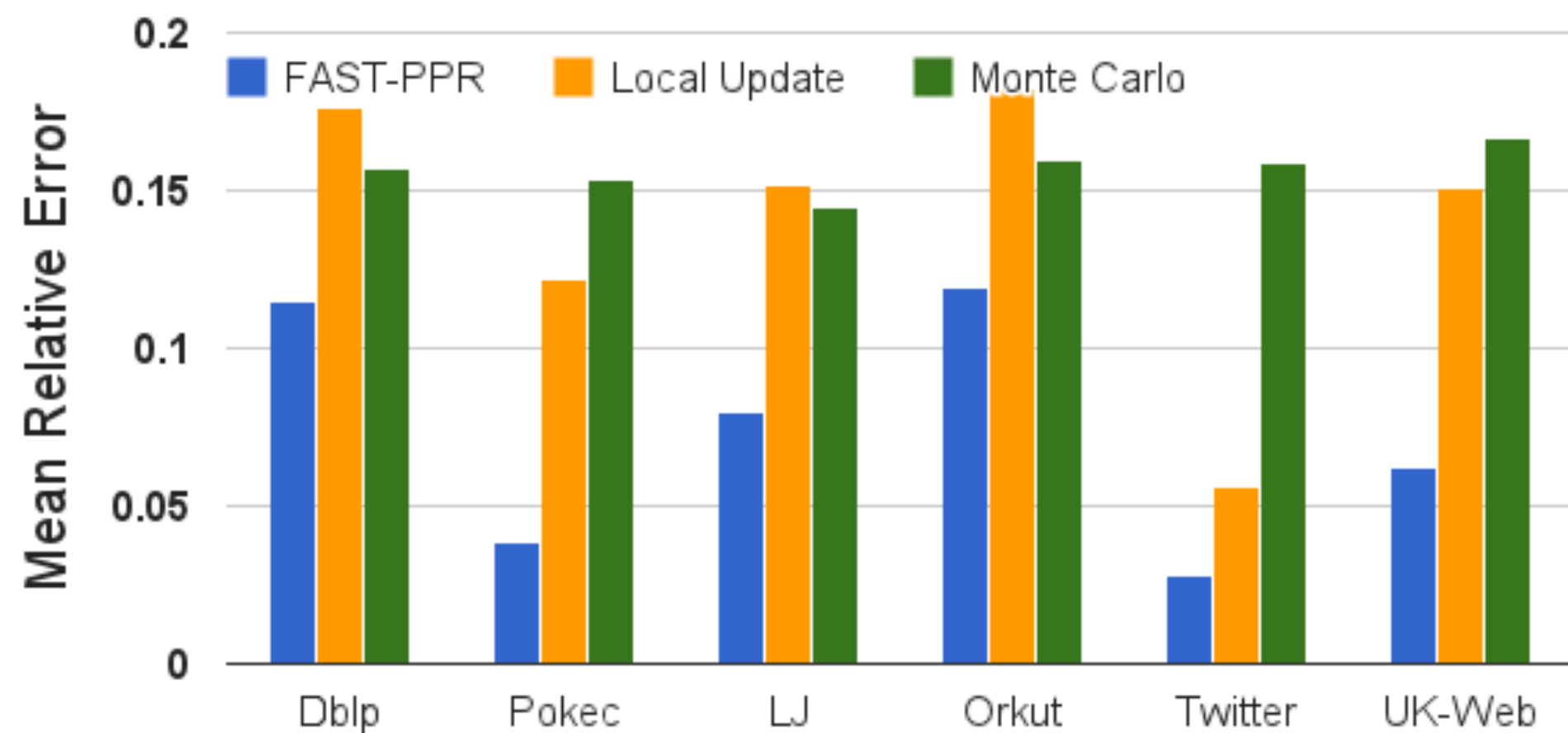
We use $\tilde{O}(\sqrt{d/\delta})$ pre-processing/space to go from average to worst case running time

Running Time (Targets sampled by PageRank)



Experiments

Relative Error of Personalized PageRank Estimates

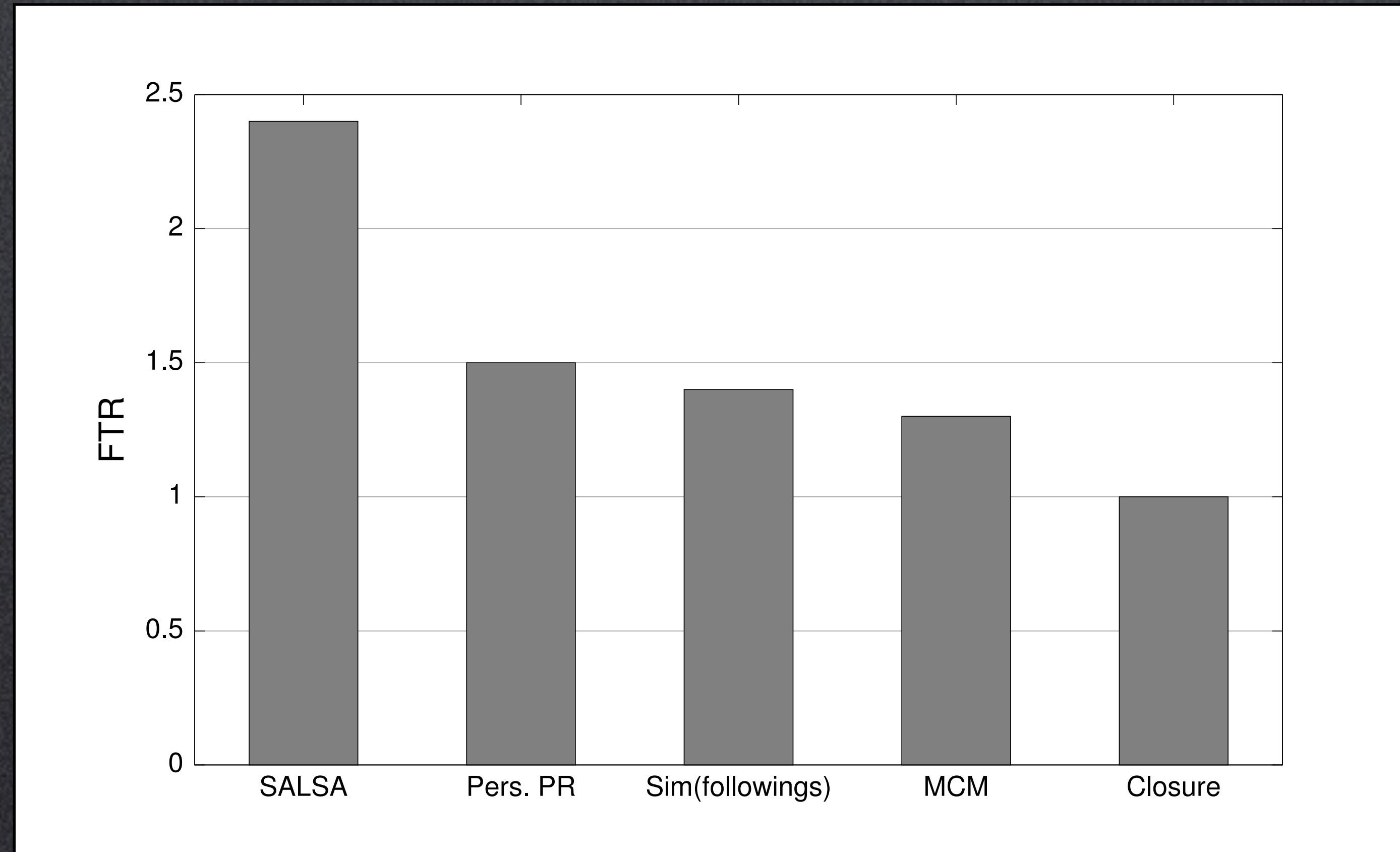


- Admits Distributed Implementation
- Works when source is a set of nodes
- Lower bound of $1/\sqrt{\delta}$
- Open problem: do we need the \sqrt{d} ?

[Lofgren, Banerjee, Goel, Seshadhri, KDD, 2014]

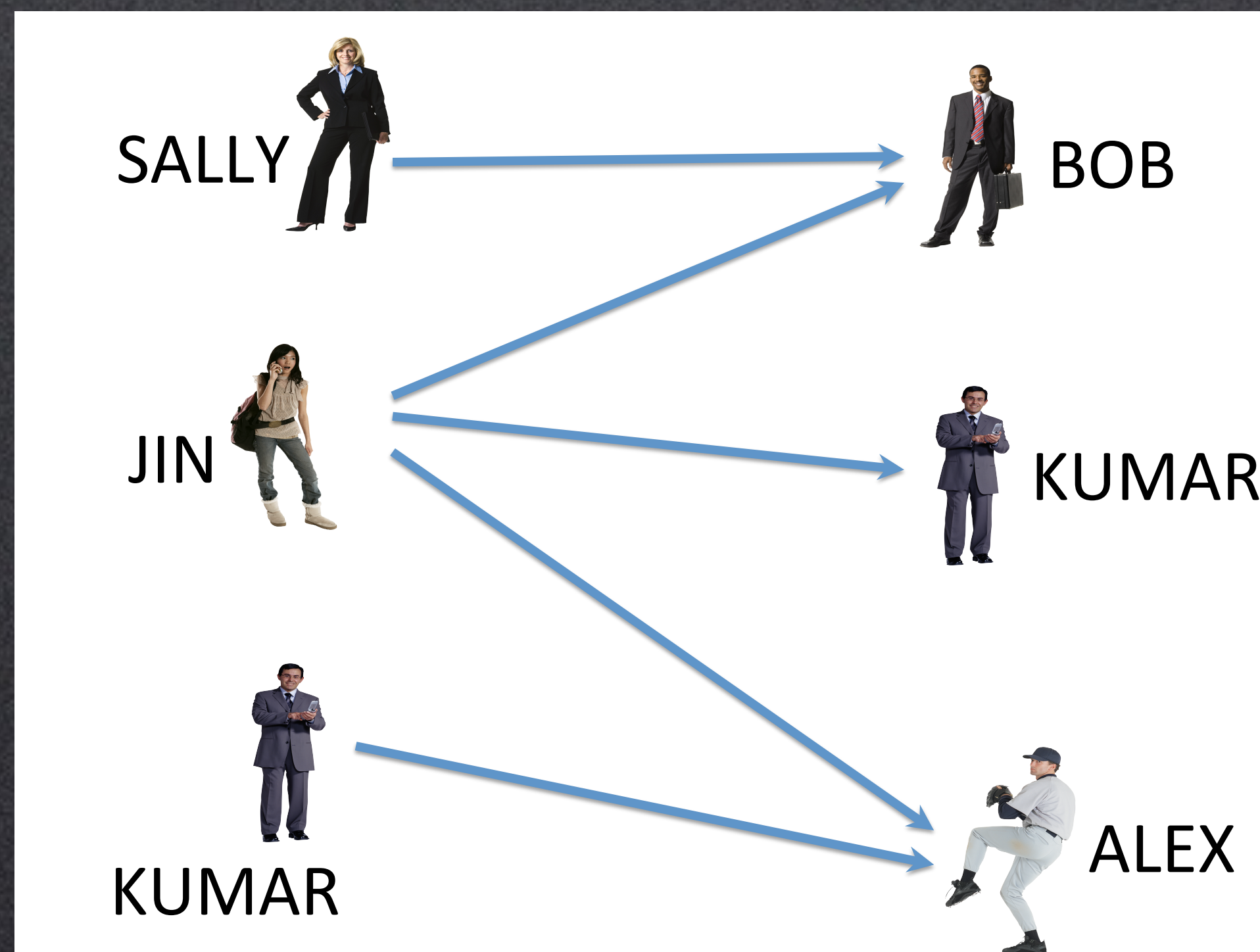
BACKUP SLIDES

Ongoing evaluation



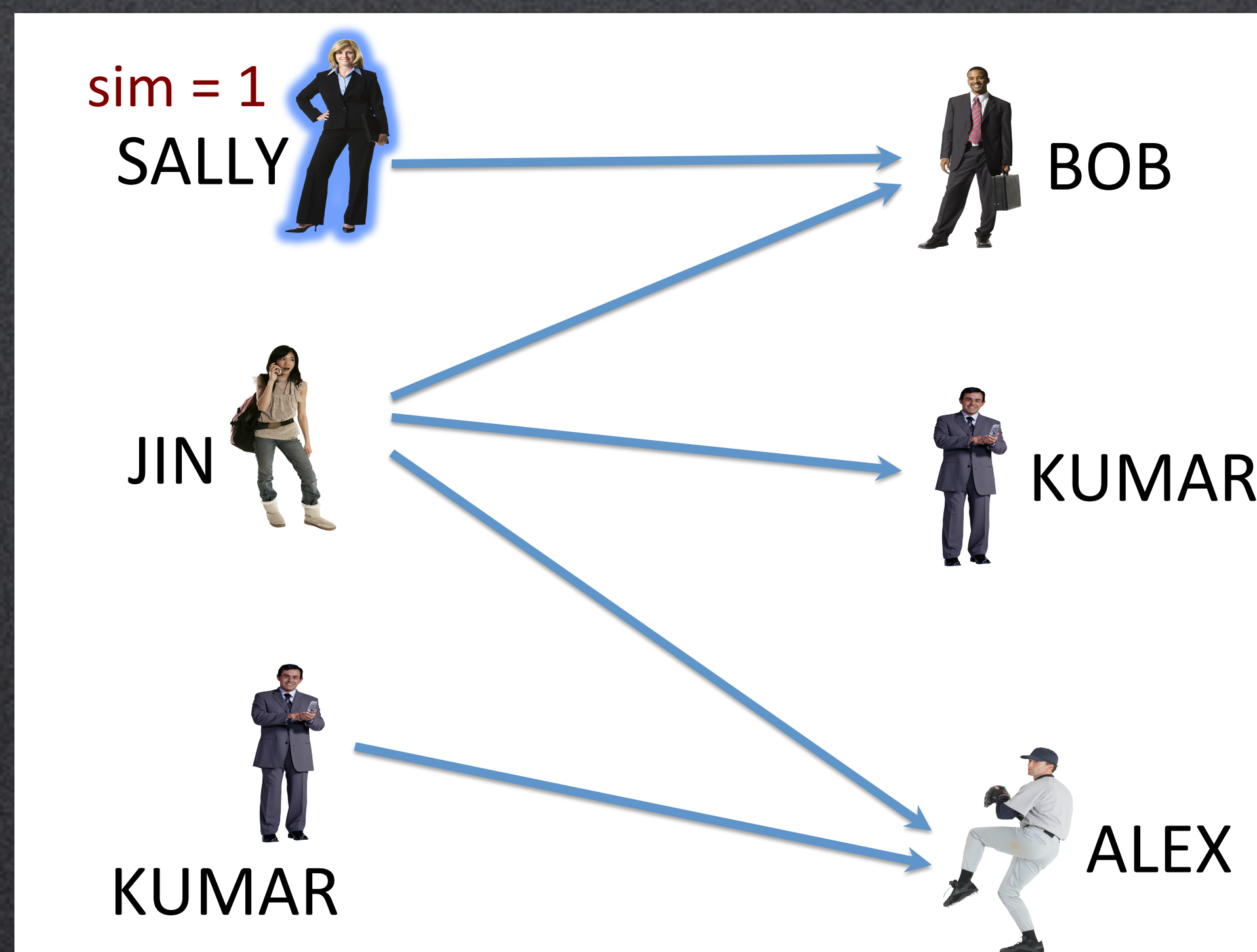
Collaborative Filter: Illustration

Use a simple propagation method: divide score by 2 and propagate (ignore the client after step 1)



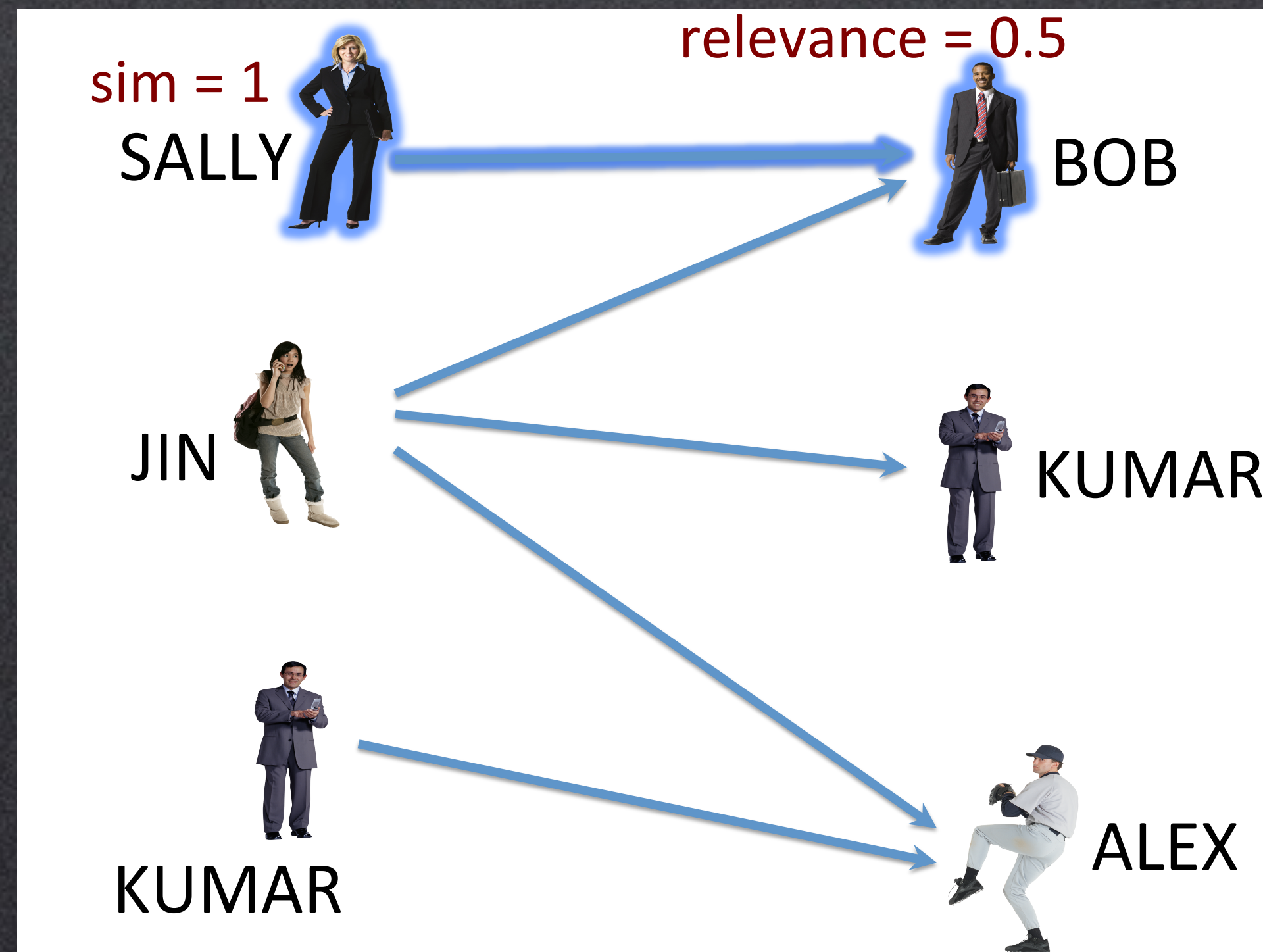
Collaborative Filter: Illustration

Use a simple propagation method: divide score by 2 and propagate (ignore the client after step 1)



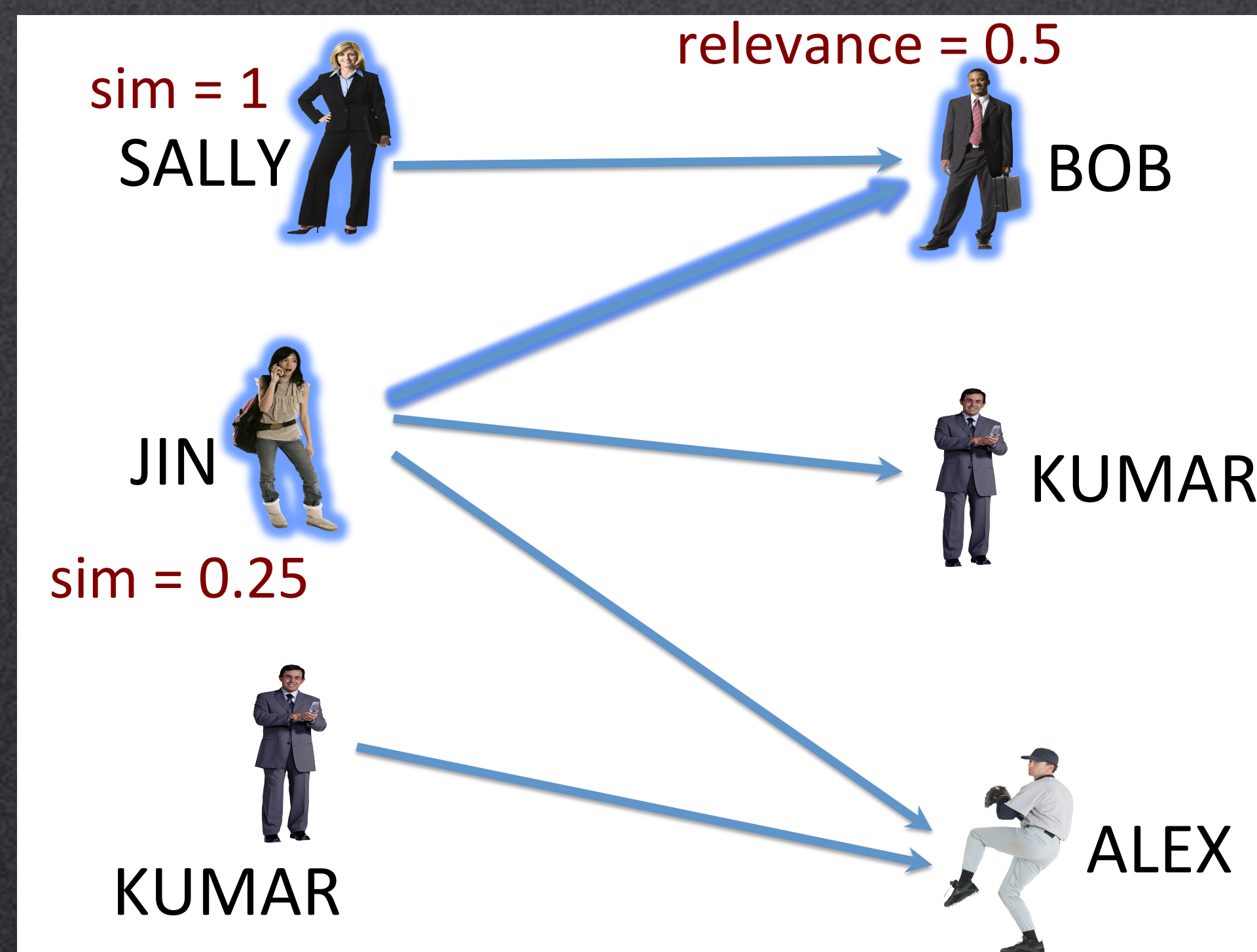
Collaborative Filter: Illustration

Use a simple propagation method: divide score by 2 and propagate (ignore the client after step 1)



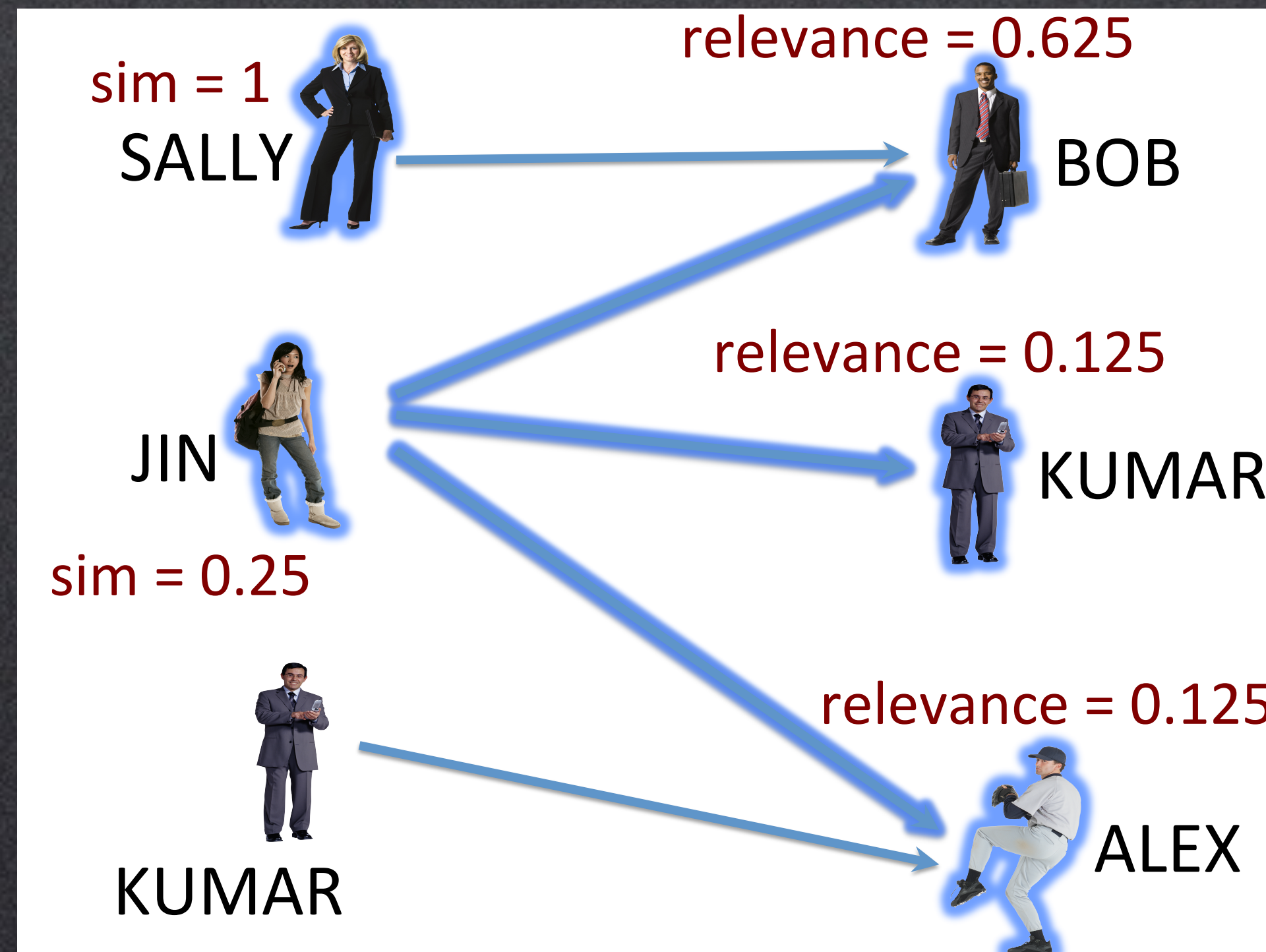
Collaborative Filter: Illustration

Use a simple propagation method: divide score by 2 and propagate (ignore the client after step 1)



Collaborative Filter: Illustration

Use a simple propagation method: divide score by 2 and propagate (ignore the client after step 1)



Collaborative Filter: Illustration

Use a simple propagation method: divide score by 2 and propagate (ignore the client after step 1)

