

# Fast Personalized PageRank On MapReduce

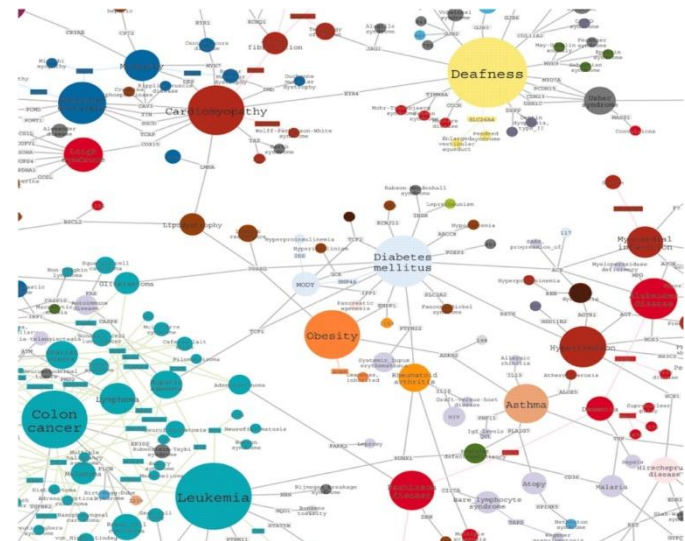
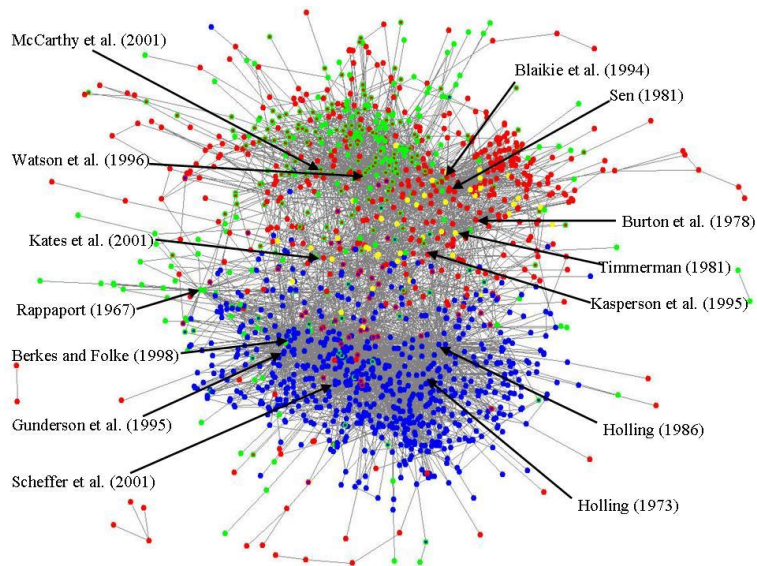
**Authors: Bahman Bahmani, Kaushik Chakrabart, Dong Xin  
In SIGMOD 2011**

**Presenter: Adams Wei Yu**

March 2015, CMU



# Graph data is Ubiquitous



**Basic Problem in Graphs:  
How do we measure the  
proximity (similarity)  
between two nodes?**

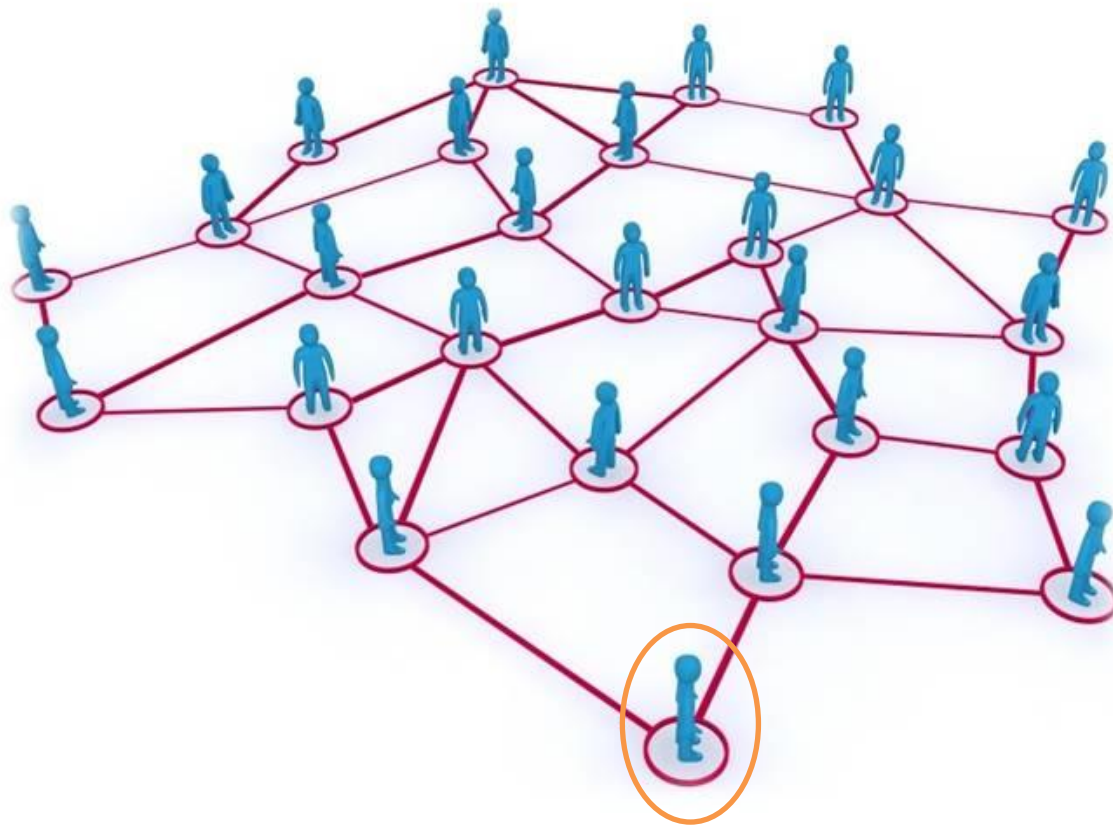
# Typical proximity measure

- Shortest distance
- Common neighbor set

# Typical proximity measure

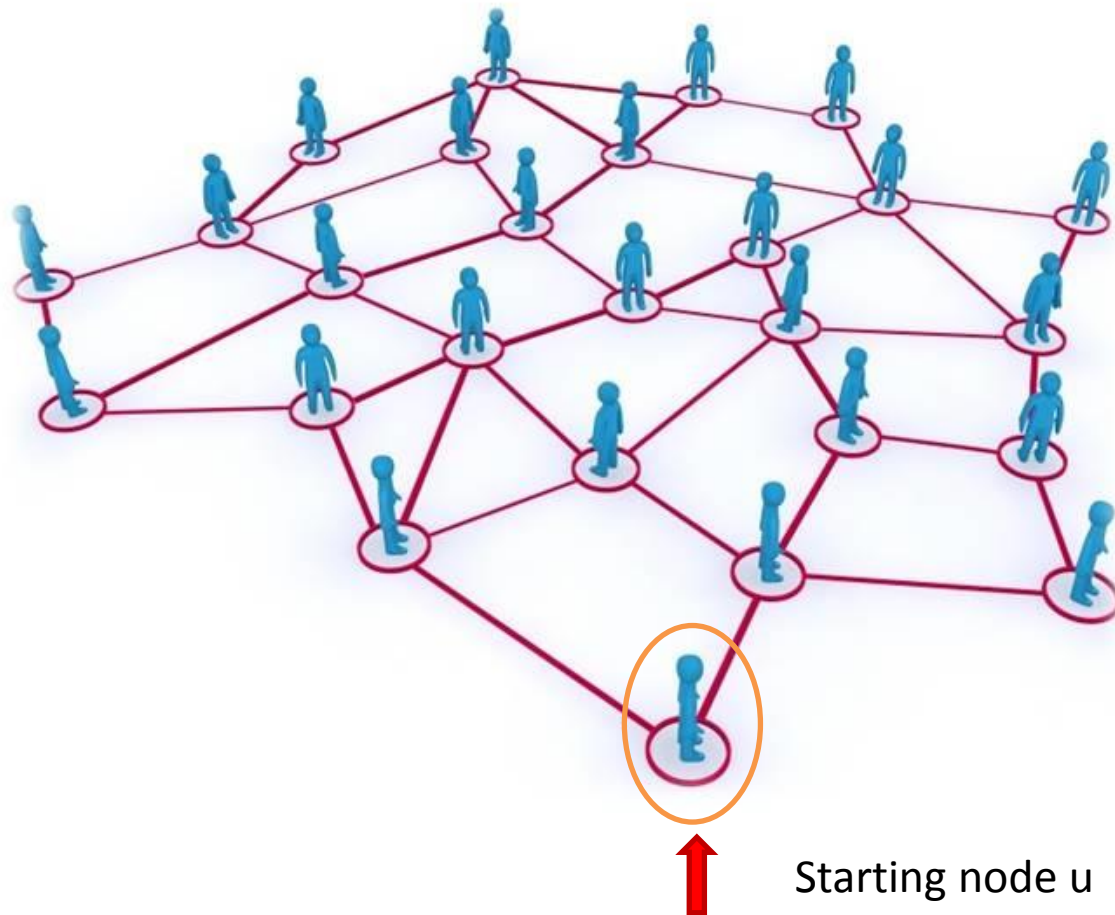
- Shortest distance
- Common neighbor set
- **Personalized PageRank (PPR), a.k.a Random Walk with Restart**
  - Considers all the possible paths from node to node
  - Captures the global structure of the graph

# Personalized PageRank



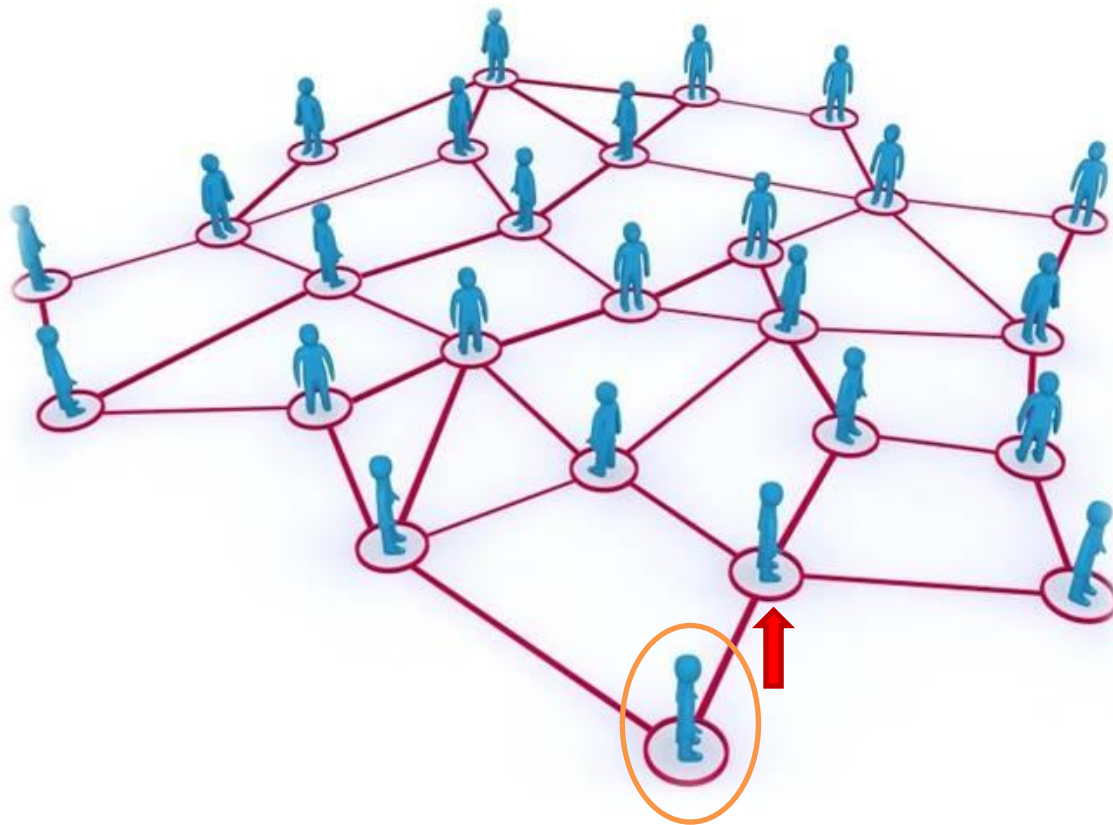
Starting node  $u$

# Personalized PageRank





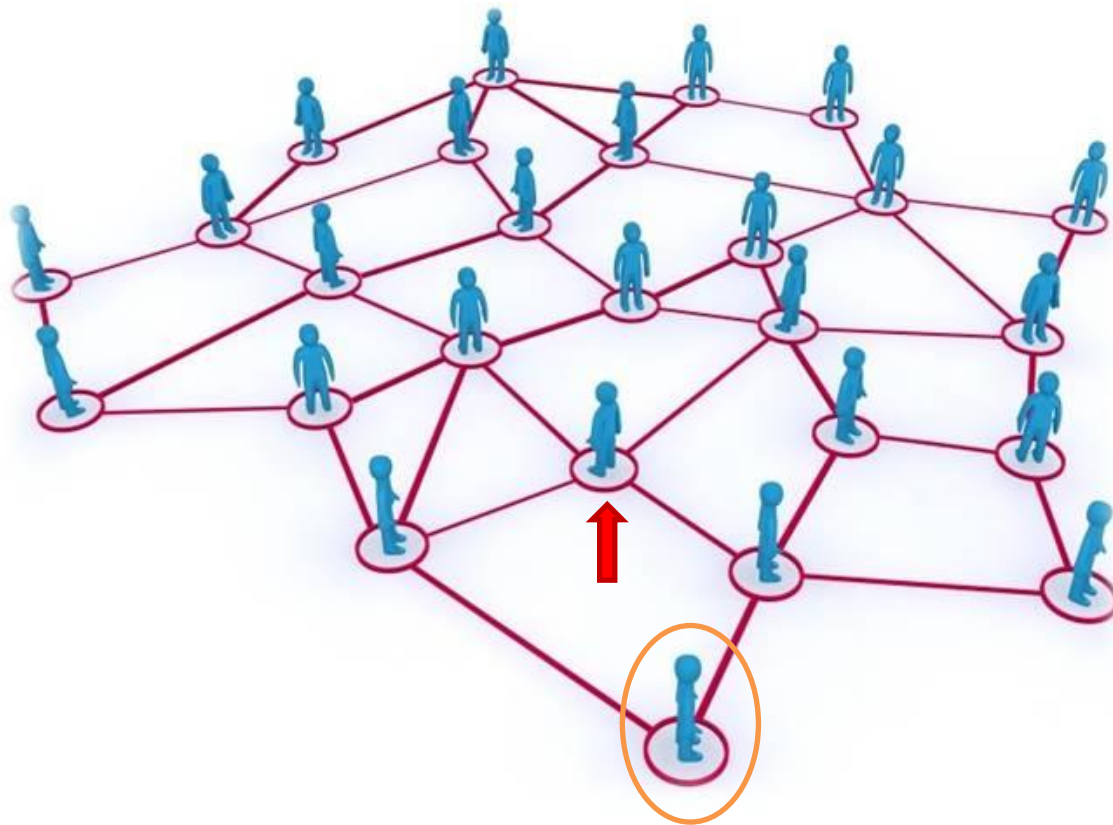
# Personalized PageRank



Starting node  $u$

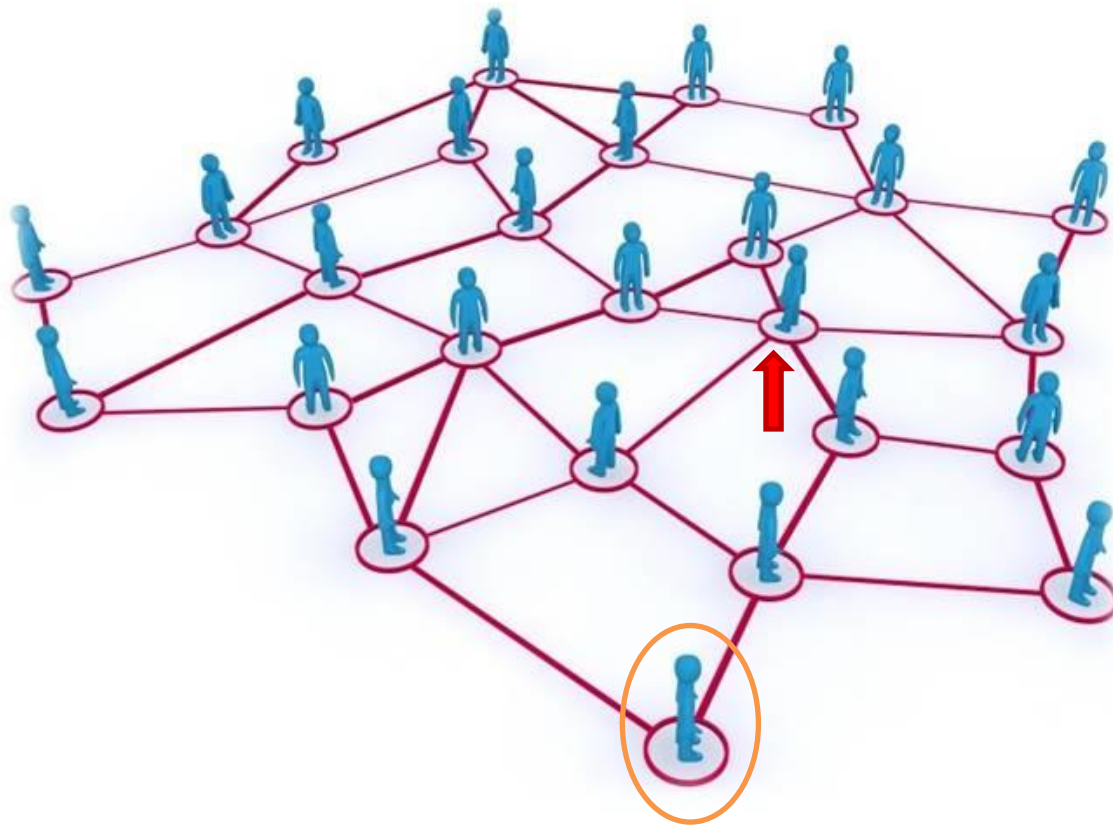


# Personalized PageRank



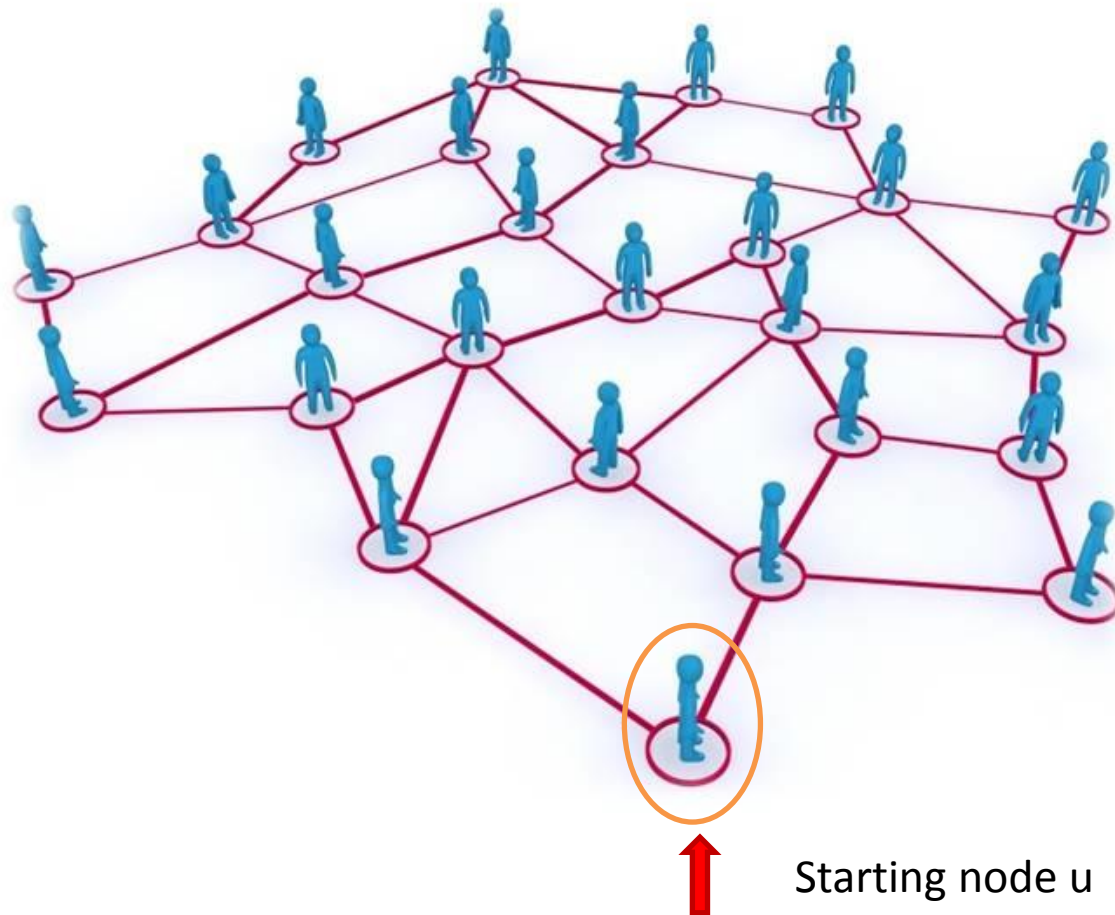
Starting node  $u$

# Personalized PageRank

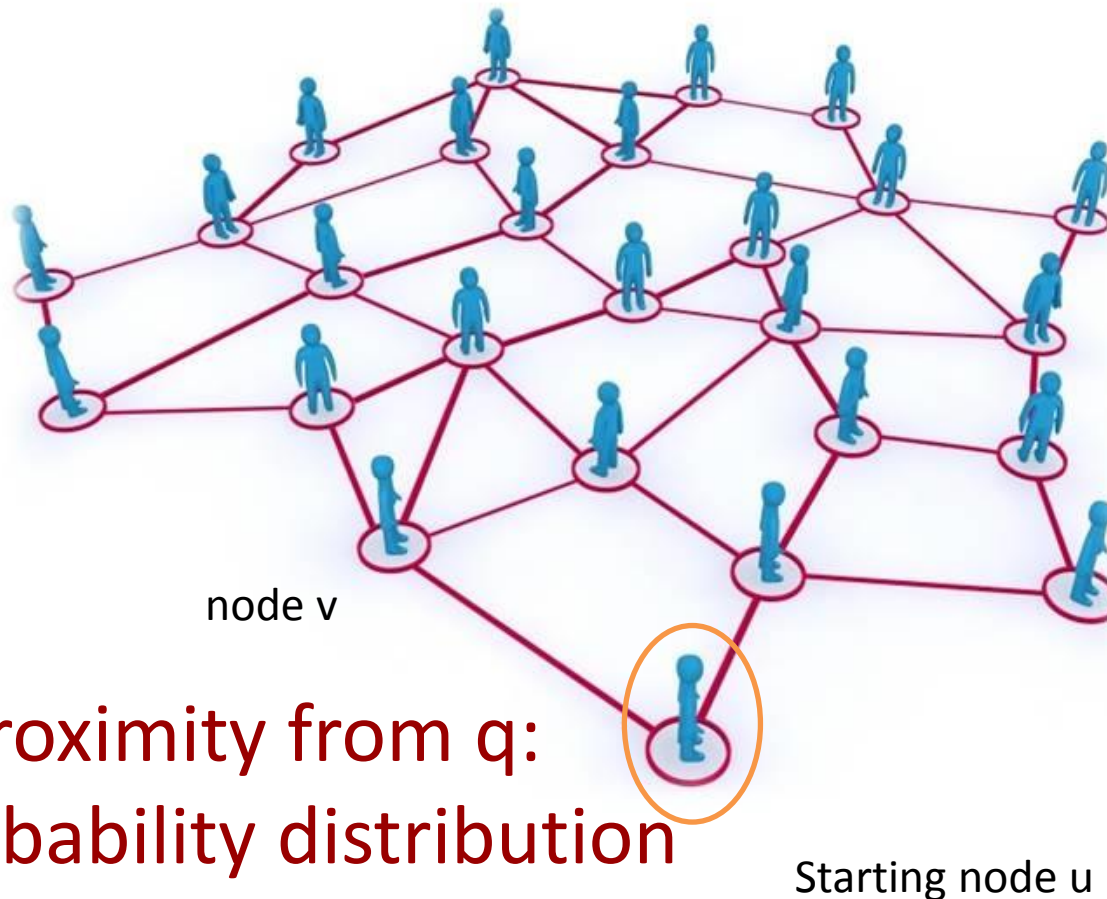


Starting node  $u$

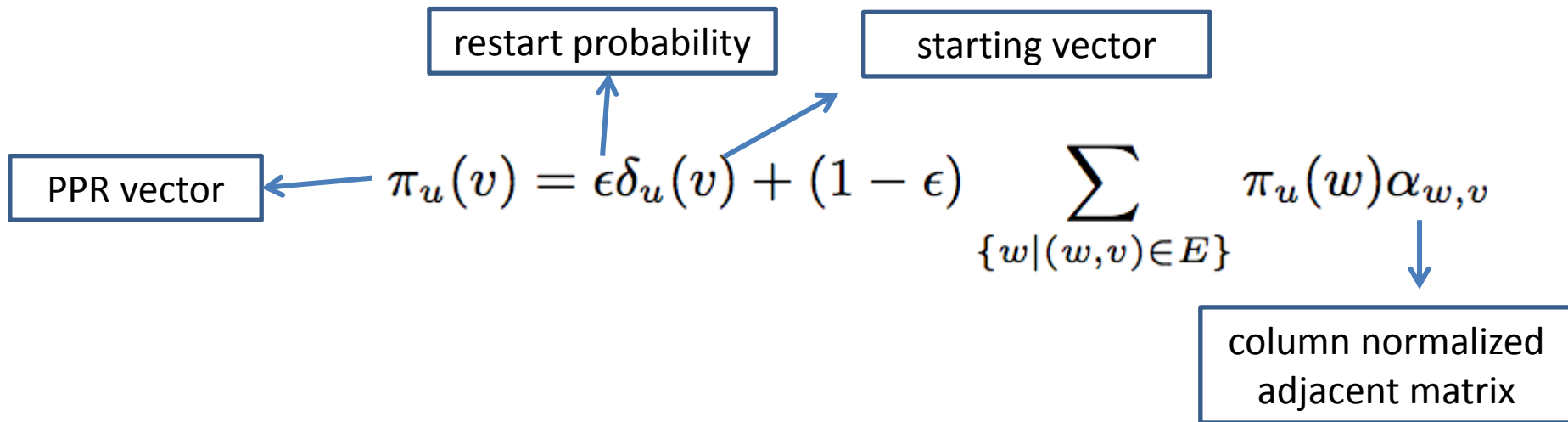
# Personalized PageRank



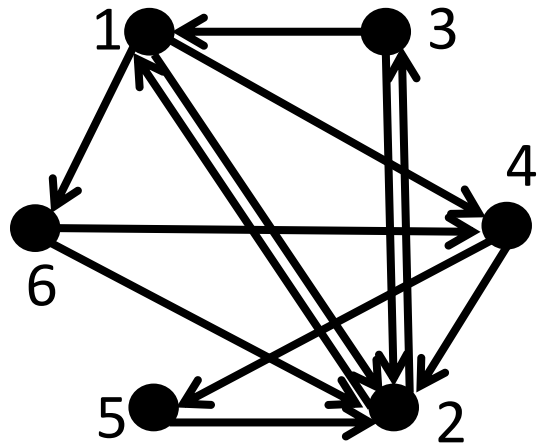
# Personalized PageRank



# Personalized PageRank



# PPR Proximity Matrix



1	2	3	4	5	6
0.32	0.24	0.24	0.19	0.20	0.18
0.28	0.39	0.29	0.31	0.33	0.30
0.12	0.17	0.27	0.13	0.14	0.13
0.13	0.10	0.10	0.23	0.08	0.14
0.06	0.04	0.04	0.10	0.18	0.06
0.09	0.07	0.07	0.05	0.06	0.20

# Application of PPR

- Citation analysis [Jeh, KDD 02]
- Link prediction [Liben Nowell, CIKM 03].
- Graph clustering [Andersen, FOCS 06].
- Recommendation system [Konstas, SIGIR 09].
- Top-k search [Fujiwara, VLDB 12]
- Reverse Top-k Search [Yu, VLDB 14]
- .....



# How to Compute the PPR?

1	2	3	4	5	6
0.32	0.24	0.24	0.19	0.20	0.18
0.28	0.39	0.29	0.31	0.33	0.30
0.12	0.17	0.27	0.13	0.14	0.13
0.13	0.10	0.10	0.23	0.08	0.14
0.06	0.04	0.04	0.10	0.18	0.06
0.09	0.07	0.07	0.05	0.06	0.20

?

# MapReduce for PPR Computation

- Reducer: Compute random neighbor of  $u$

# MapReduce for PPR Computation

- Reducer: Compute random neighbor of  $u$ 
  - $G = \langle u, v, \text{weight} \rangle$ . Reducer uses  $u$  as the key, groups all  $\langle u, v, \text{weight} \rangle$  triples (for all  $v$ ), and generating one random neighbor  $v$  for  $u$  and outputs  $\langle u, v \rangle$ .

# MapReduce for PPR Computation

- Reducer: Compute random neighbor of  $u$ 
  - $G = \langle u, v, \text{weight} \rangle$ . Reducer uses  $u$  as the key, groups all  $\langle u, v, \text{weight} \rangle$  triples (for all  $v$ ), and generating one random neighbor  $v$  for  $u$  and outputs  $\langle u, v \rangle$ .
- Combiner: Extend  $N = \langle u, v \rangle$  to  $\langle u, v, t \rangle$ .

# MapReduce for PPR Computation

- Reducer: Compute random neighbor of  $u$ 
  - $G = \langle u, v, \text{weight} \rangle$ . Reducer uses  $u$  as the key, groups all  $\langle u, v, \text{weight} \rangle$  triples (for all  $v$ ), and generating one random neighbor  $v$  for  $u$  and outputs  $\langle u, v \rangle$ .
- Combiner: Extend  $N = \langle u, v \rangle$  to  $\langle u, v, t \rangle$ .
  - Find a random neighbor  $t$  of  $v$  by joining  $N$  and  $G$  on condition  $N.v = G.u$  from the graph  $G$  and output  $\langle u, v, t \rangle$ .

# Computation of PPR

- Power Iteration:

$$\pi_u^{(i)}(v) = \epsilon \delta_u(v) + (1 - \epsilon) \sum_{\{w | (w,v) \in E\}} \pi_u^{(i-1)}(w) \alpha_{w,v}$$

# Computation of PPR

- Power Iteration:

$$\pi_u^{(i)}(v) = \epsilon \delta_u(v) + (1 - \epsilon) \sum_{\{w | (w,v) \in E\}} \pi_u^{(i-1)}(w) \alpha_{w,v}$$

- Basic MapReduce:

- Given a graph  $G = \langle u, v, \alpha_{u,v} \rangle$ , the initialization of  $\overrightarrow{\pi}_u^{(0)}$  is a Reducer of graph on key  $u$ .
- Update of  $\overrightarrow{\pi}_u^{(i)}$  can be implemented as combiner joining  $\overrightarrow{\pi}_u^{(i-1)}$  and  $G$ .



# Computation of PPR

- Power Iteration:

$$\pi_u^{(i)}(v) = \epsilon \delta_u(v) + (1 - \epsilon) \sum_{\{w | (w,v) \in E\}} \pi_u^{(i-1)}(w) \alpha_{w,v}$$

- Basic MapReduce:

- Given a graph  $G = \langle u, v, \alpha_{u,v} \rangle$ , the initialization of  $\overrightarrow{\pi_u^{(0)}}$  is a Reducer of graph on key  $u$ .
- Update of  $\overrightarrow{\pi_u^{(i)}}$  can be implemented as combiner joining  $\overrightarrow{\pi_u^{(i-1)}}$  and  $G$ .

**Space:  $O(n^2)$ !**

**We may not need the exact  
PPR for most tasks!**

# Computation of PPR

- Monte Carlo Simulation

Simulate  $R$  random walks starting from  $u$ ,  
the portion of visits to  $v$  is approximately  $\pi_u(v)$

# Computation of PPR

- Monte Carlo Simulation

Simulate  $R$  random walks starting from  $u$ ,  
the portion of visits to  $v$  is approximately  $\pi_u(v)$

- Basic MapReduce:

- A Reducer to initialize  $R$  random walks from  $u$ .
- A sequence of Combiner iterations to extend each random walk until it restarts at  $u$ .
- A final Reducer to aggregate the frequencies of visits to every node  $v$  in all  $R$  walks (for each source node), and approximate the PPR values.

# Computation of PPR

- Monte Carlo Simulation

Simulate  $R$  random walks starting from  $u$ ,  
the portion of visits to  $v$  is approximately  $\pi_u(v)$

- Basic MapReduce:

- A Reducer to initialize  $R$  random walks from  $u$ .
- A sequence of Combiner iterations to extend each random walk until it restarts at  $u$ .
- A final Reducer to aggregate the frequencies of visits to every node  $v$  in all  $R$  walks (for each source node), and approximate the PPR values.

**Needs to execute the Combiner many times for long walks!**

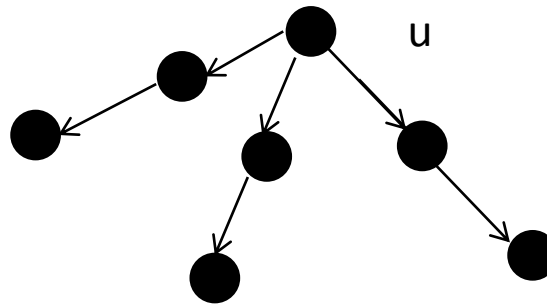
# Computation of PPR

- Monte Carlo Simulation

Given a graph  $G$  and a length  $\lambda$ , outputs one random walk of length  $\lambda$  starting from each node in the graph.

# Basic Idea

- Preprocessing: Sample  $\eta = \lambda/\theta$  short segments of length  $\theta$  out of each node.
- Online Computation: Merge these segments to form the longer walk.

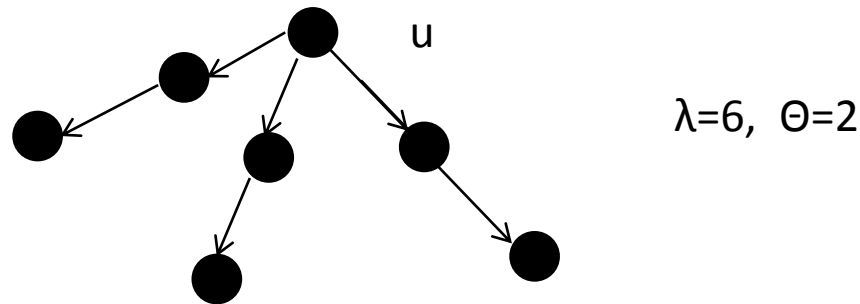


$\lambda=6, \theta=2$



# Basic Idea

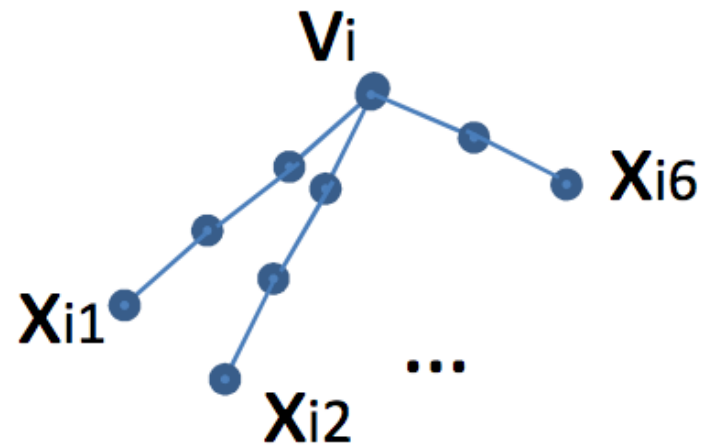
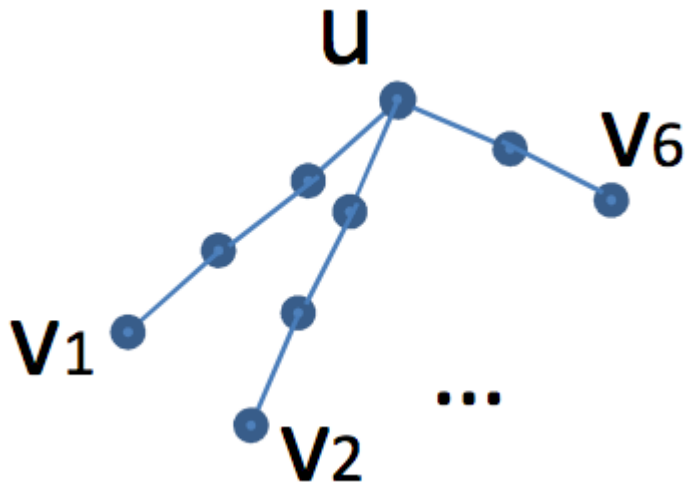
- Preprocessing: Sample  $\eta = \lambda/\theta$  short segments of length  $\theta$  out of each node.
- Online Computation: Merge these segments to form the longer walk.



1. How to ensure it is a REAL random walk?
2. How to make it efficient?

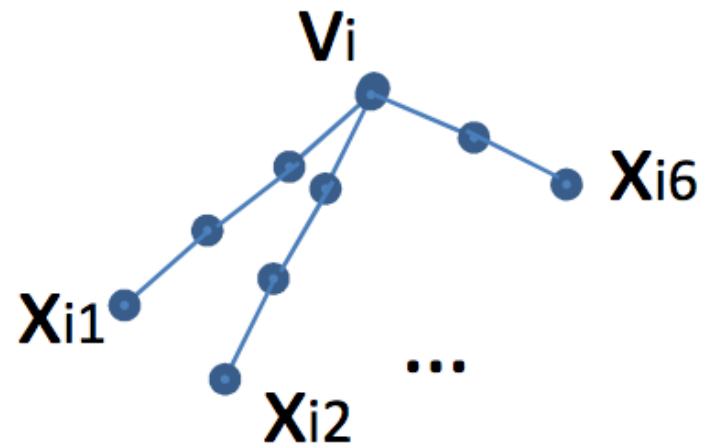
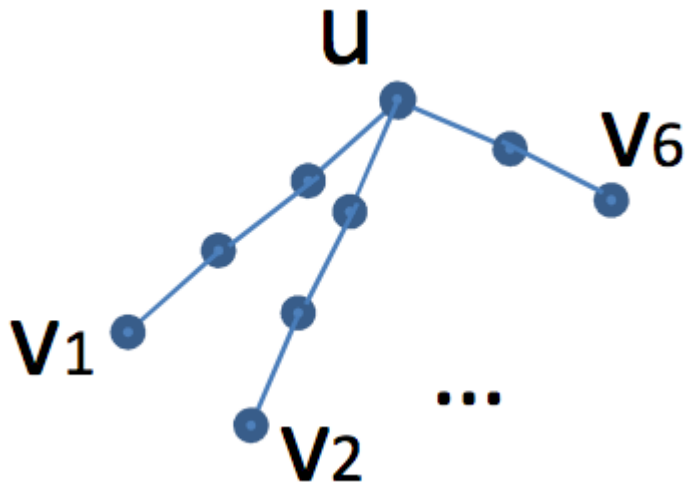
# Doubling Algorithm: Initialization

- Assume  $\lambda = 17, \theta = 3$ .
- For each node  $u$ , generate  $\eta = \lceil 17/3 \rceil = 6$  segments  $S[u, i]$  ( $1 \leq i \leq 6$ ),  $i$  is its ID.
- $S[u, 6]$  is of length 2 while the other segments  $S[u, i]$  ( $i < 6$ ) are of length 3.



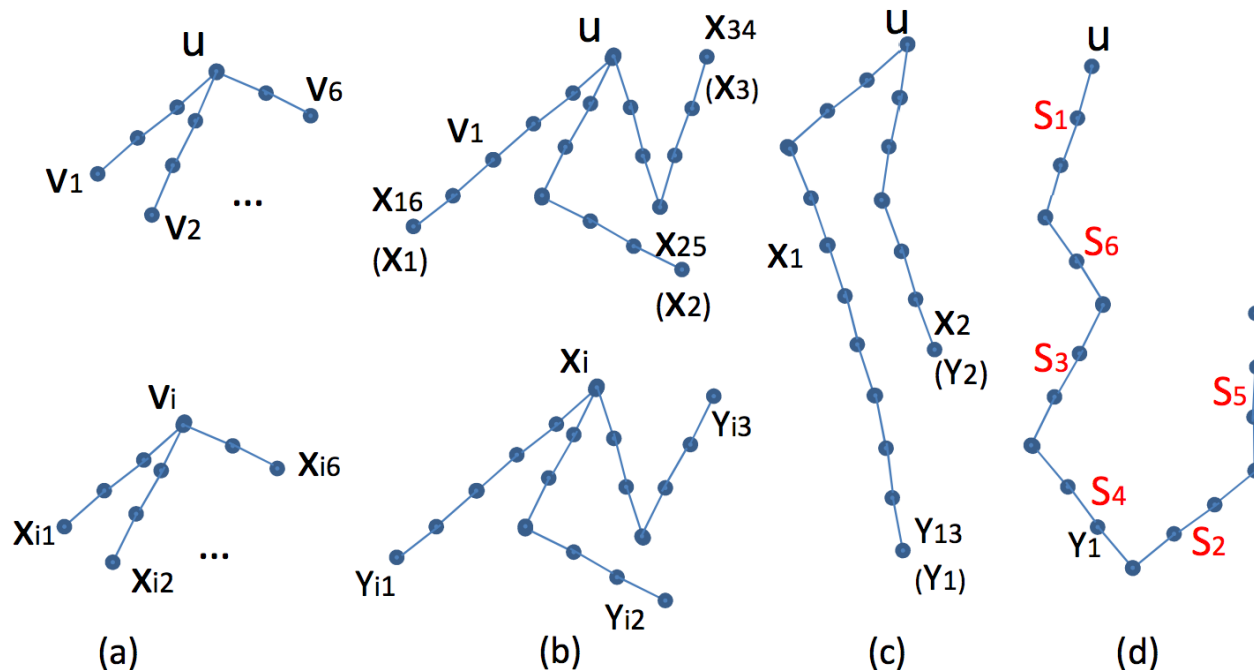
# Doubling Algorithm: Initialization

- $S[u, i]$ :  $i$ -th segment of  $u$
- $W[u, i, \eta]$ : RW from  $u$ ,  $\eta$  is the maximum ID at the current iteration.
- In the beginning,  $\eta = 6$ , and  $W[u, i, \eta] = S[u, i]$  for  $i = 1, \dots, 6$ .



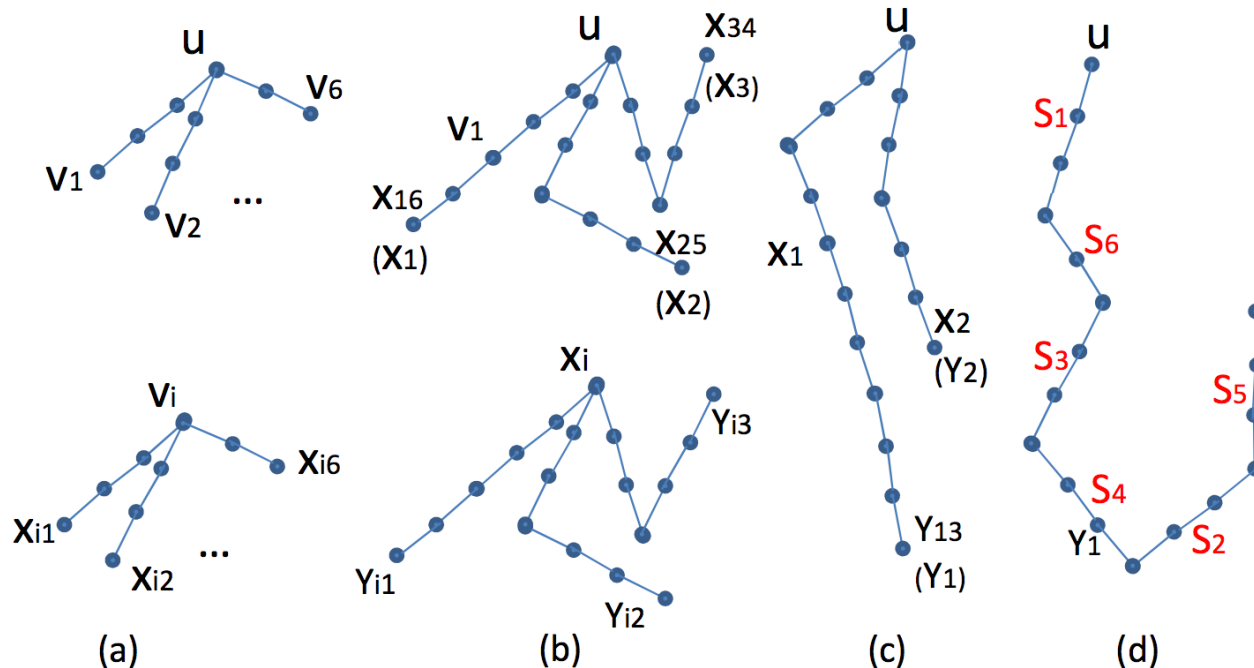
# Doubling Algorithm: Merge

- Appends  $W2$  to  $W1$  if:
  - $W1.LastNode = W2.FirstNode$
  - $W1.ID < W2.ID$
  - $W1.ID + W2.ID = \eta + 1$ . (ensure each segment is a proper random walk)



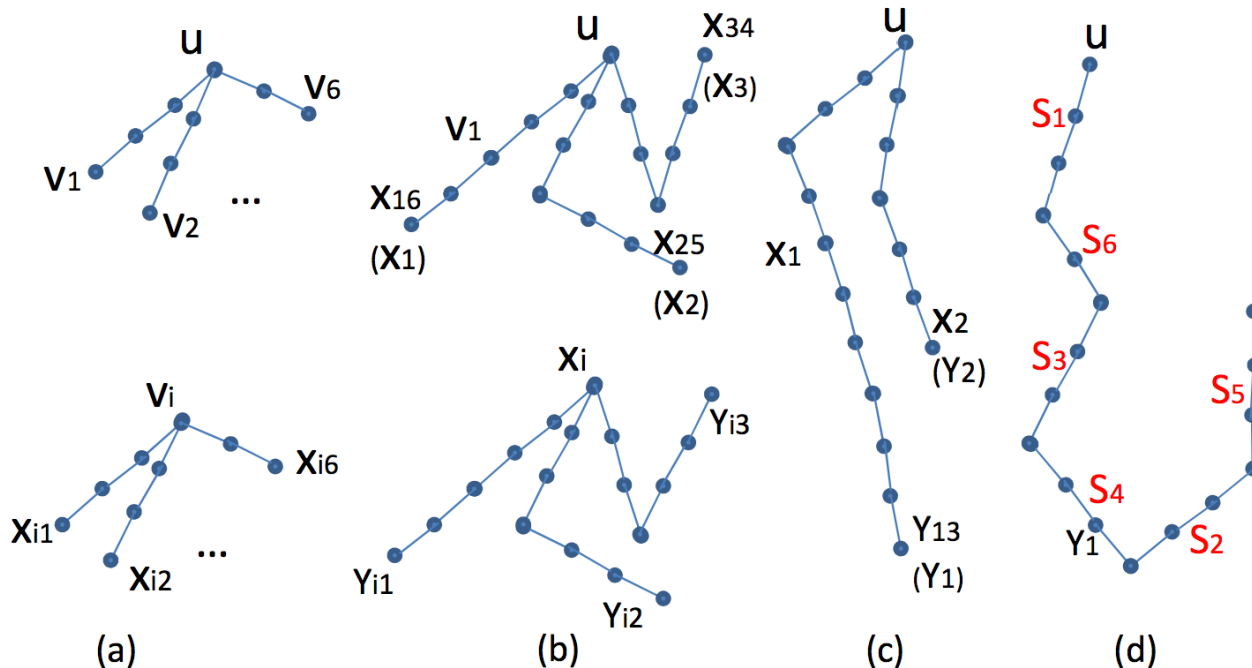
# Doubling Algorithm: Merge

- For a node  $u$ , we merge  $W[u,i,6]$  with  $W[v_i,7-i,6]$  for  $i = 1,2,3$ , and get 3 new segments:
  - $W[u,1,3]$  that ends at  $x_1$ ,  $W[u,2,3]$  that ends at  $x_2$ , and  $W[u,3,3]$  that ends at  $x_3$



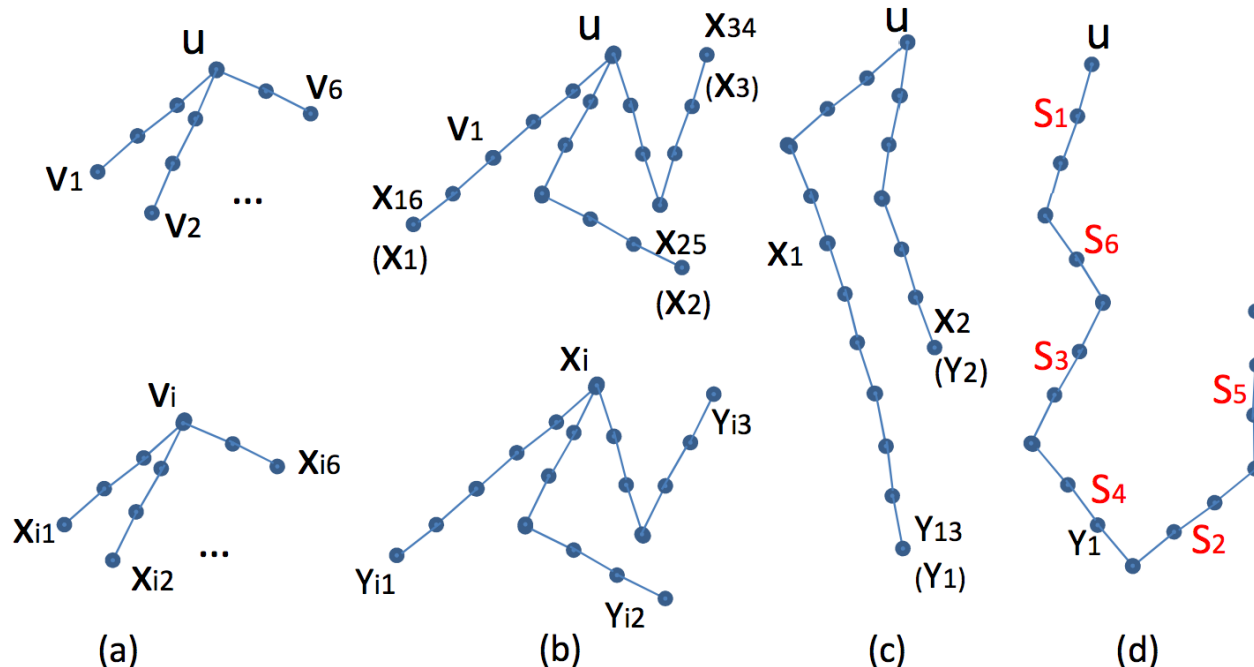
# Doubling Algorithm: Merge

- Continue:  $\eta = 3$ , and we will merge  $W[u,1,3]$  with  $W[x_1,3,3]$ . See Fig (c).
- Finally: merge  $W[u,1,2]$  with  $W[y_1,2,2]$ , and get  $W[u,1,1]$ , which has length  $\lambda = 17$ .



# Doubling Algorithm: Merge

Shrinkage fast: One merging iteration will reduce the maximum ID from  $\eta$  to  $\lfloor \frac{\eta+1}{2} \rfloor$  !





# Optimality of the Algorithm

- Theorem: The Doubling algorithm with parameters  $\lambda, \theta$  finishes in  $\theta + \lceil \log_2 \lceil \frac{\lambda}{\theta} \rceil \rceil$  MapReduce iterations.

$$1 + \lceil \log_2 \lambda \rceil \quad \text{If } \theta = 1$$

**Optimal for in the class of Natural Algorithms(NA)!**

**NA only allows Extend and Merge operations.**

# Computing the PPR

- Choose different  $\lambda$ ,  $\theta$  and  $R$  for different nodes.
- Repeat Doubling algorithm.
- Collect and count number of visits.
  - $C(u,v)_+ = \text{number of visits to } v \text{ in } W_i[u]$

# I/O cost of the Algorithm

$$O(m \max_{1 \leq i \leq R} \{\theta_i\} + n \sum_{i=1}^R (\lambda_i \theta_i + \lambda_i \log_2 \frac{\lambda_i}{\theta_i}))$$

# Experiments

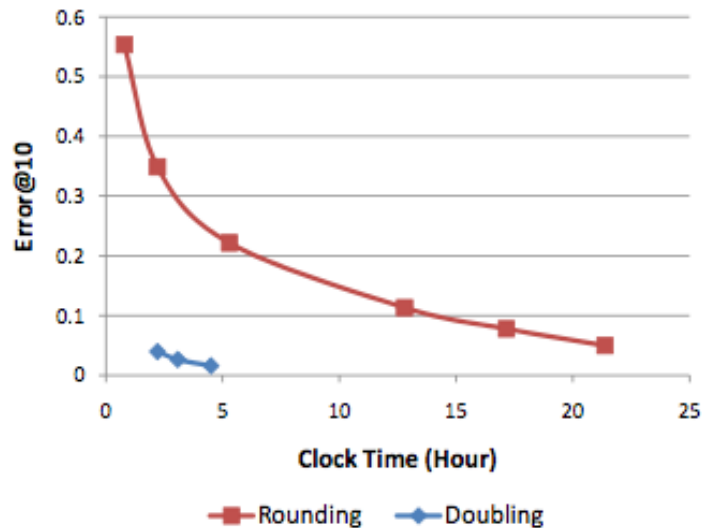


Figure 2:  $\overline{Err}(\cdot, 10)$  vs. Clock Time

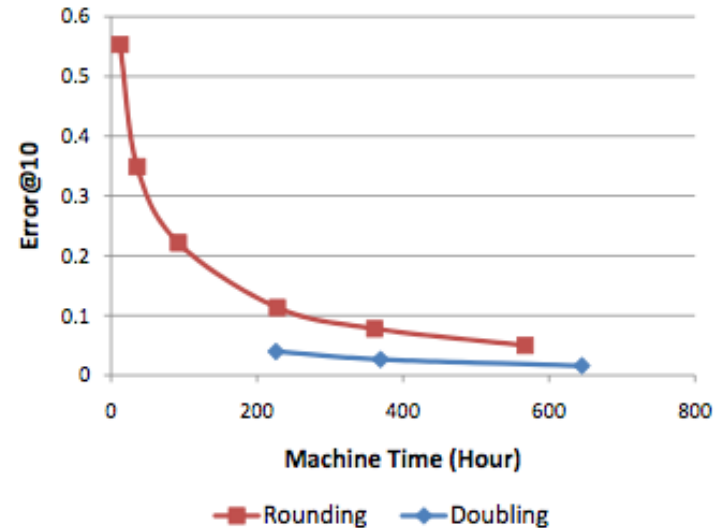


Figure 3:  $\overline{Err}(\cdot, 10)$  vs. Machine Time

Graph Size: 112M nodes and 513M edges

# Take-home Message

- A new PPR computation method on MapReduce is proposed.
- Theoretically sound
  - Optimal in the class of natural algorithms.
  - Manageable I/O cost.
- Empirically works well.

Thanks!