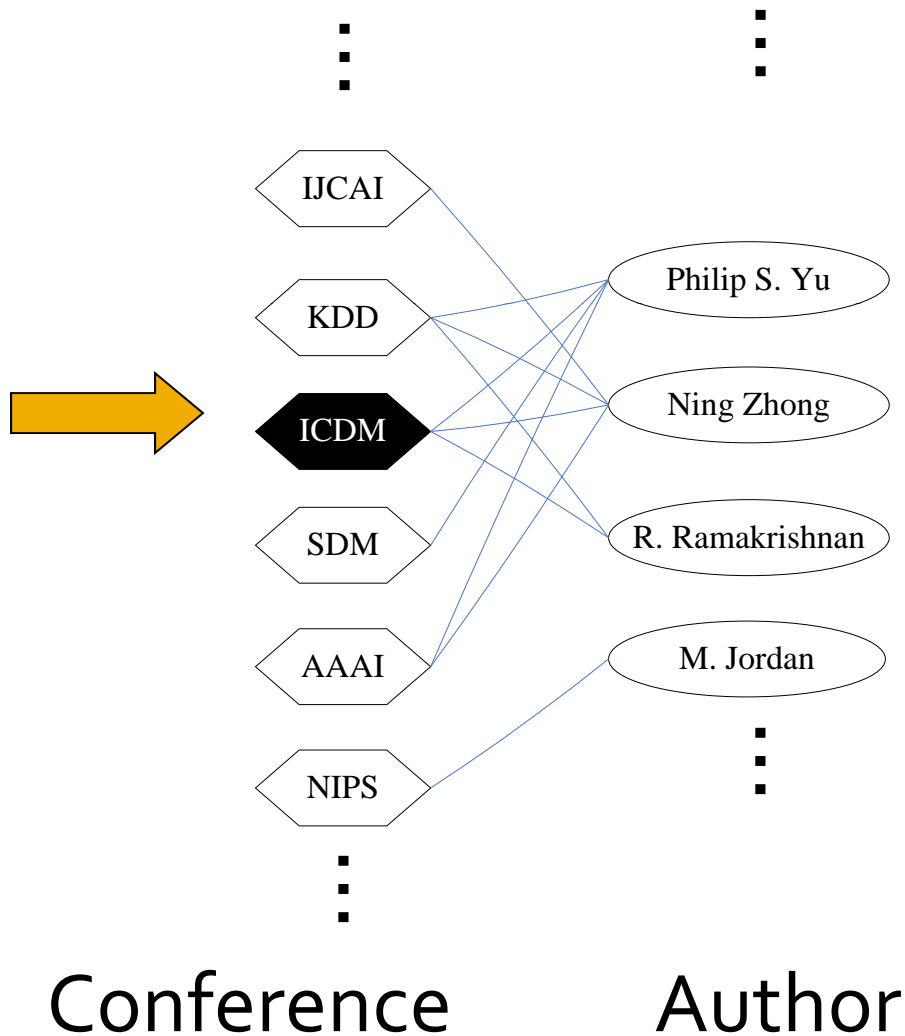


Proximity in Graphs by using Random Walks

CS345a: Data Mining
Jure Leskovec and Anand Rajaraman
Stanford University



Neighborhood Formulation

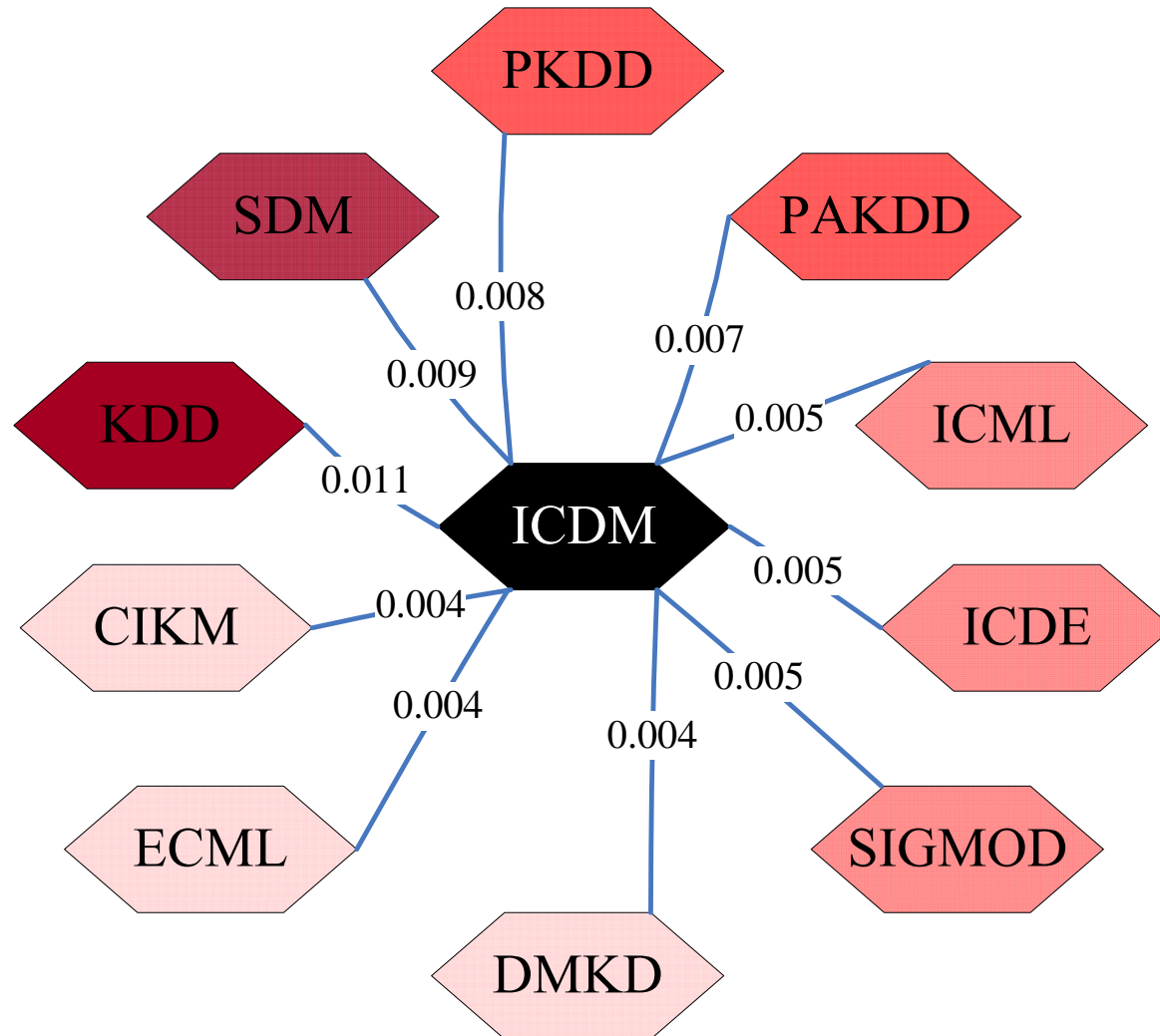


Q: what is most related conference to ICDM

A: Random walks!

[Sun, ICDM2005]

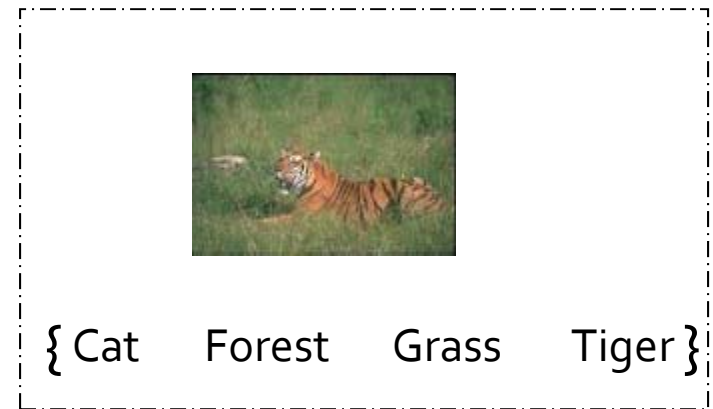
Neighborhood Search on Graphs



Automatic Image Captioning



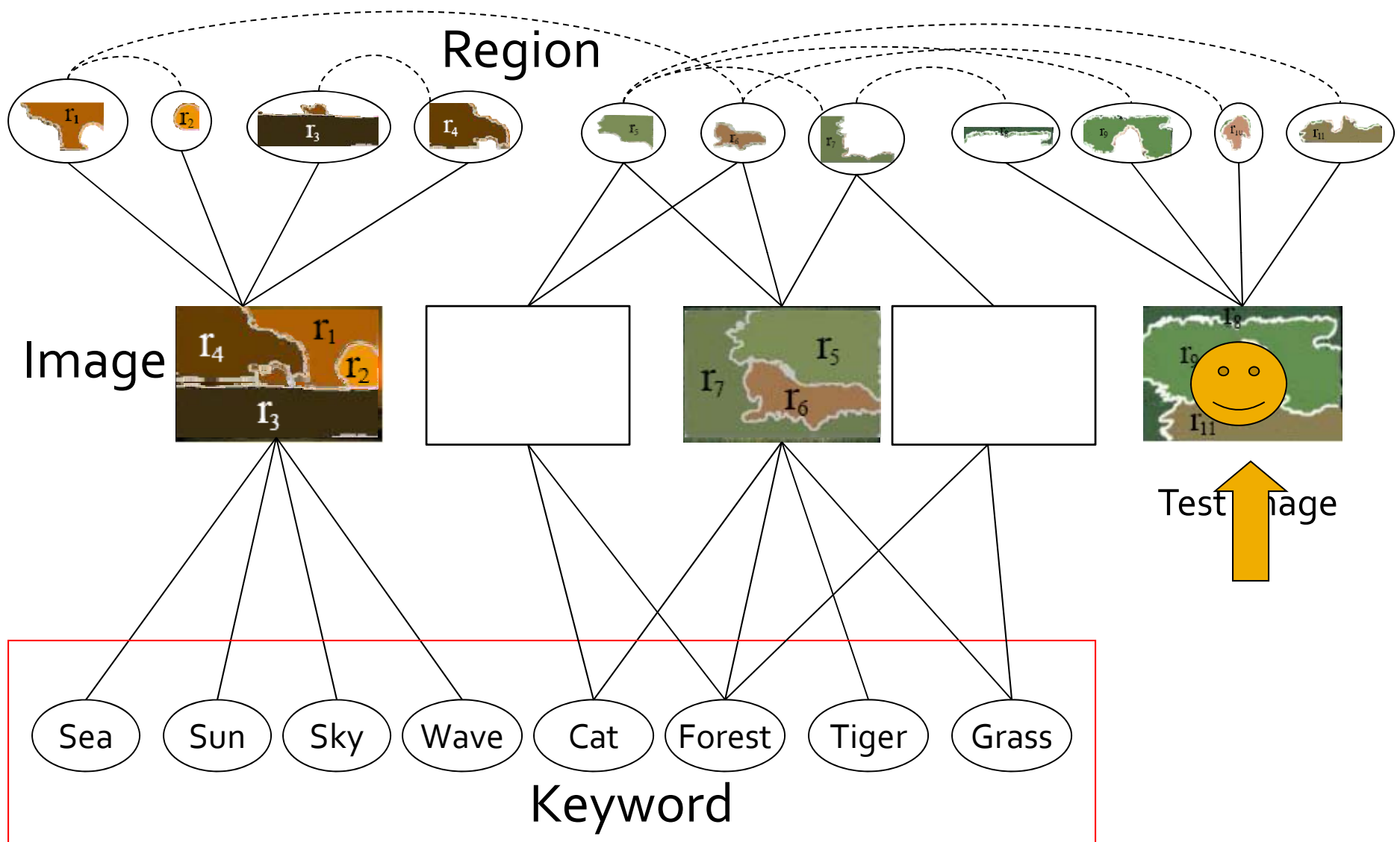
...



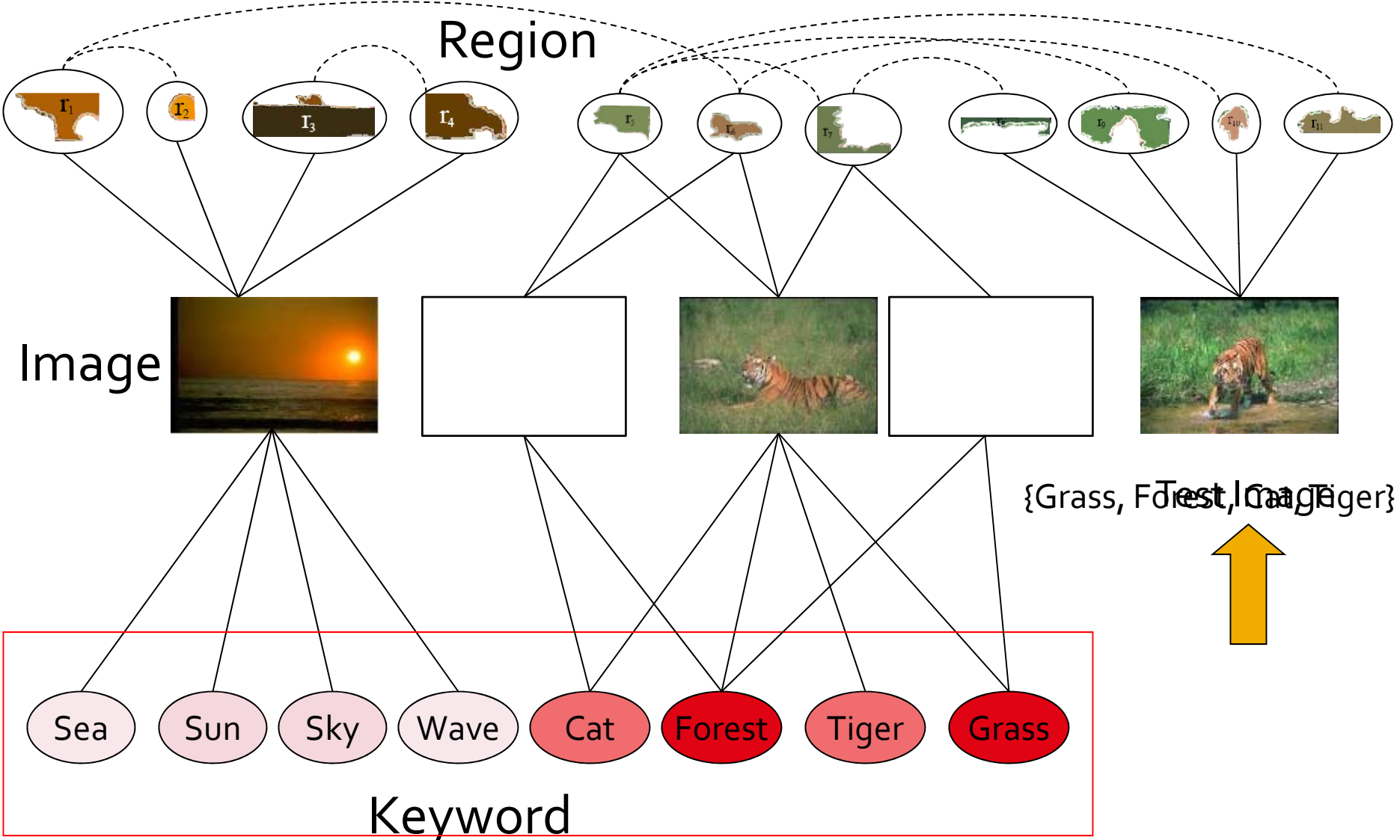
A: Proximity on graphs!

[Pan KDD2004]

Automatic Image Captioning

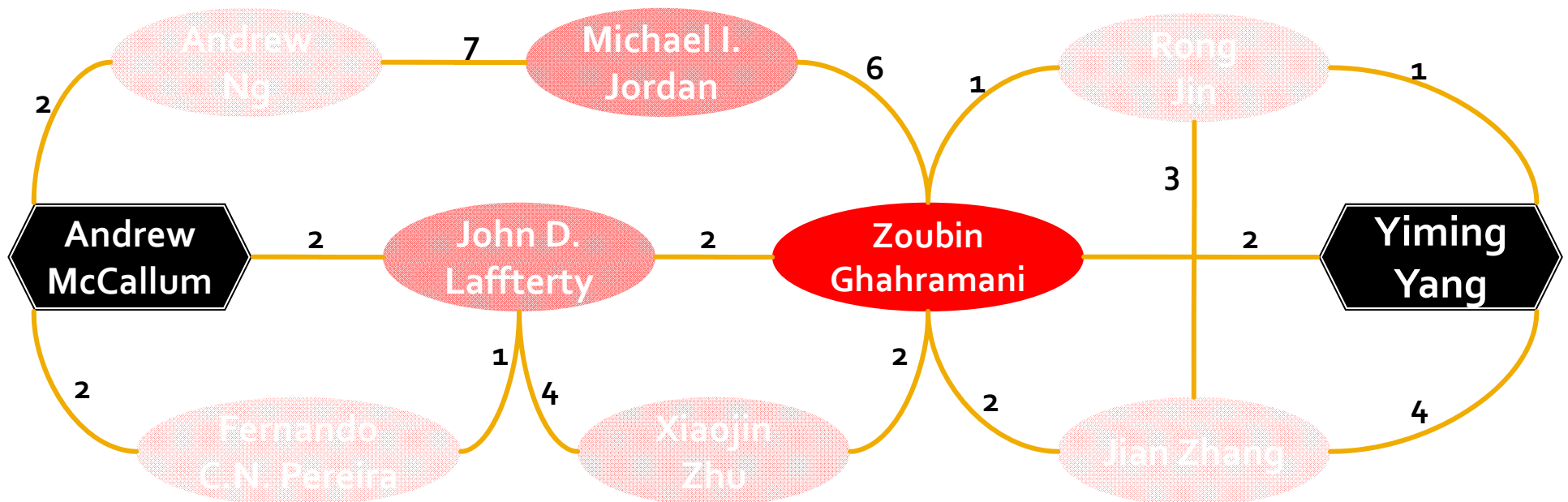


Automatic Image Captioning



Center-Piece Subgraphs

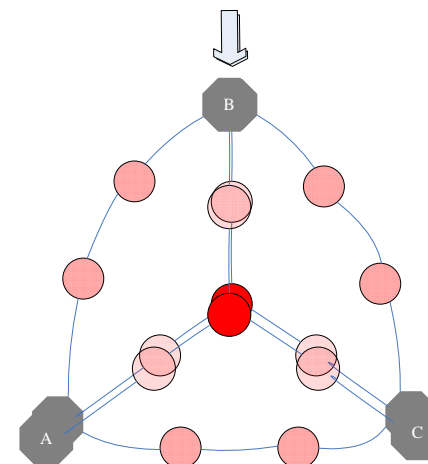
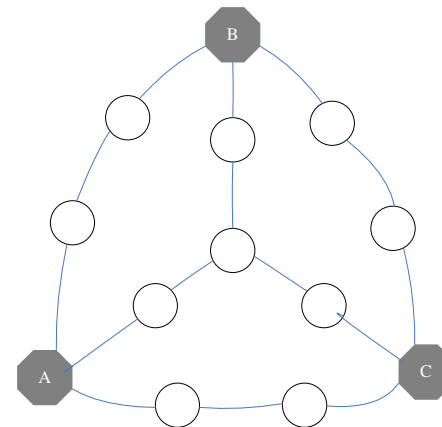
- Connection subgraphs:
 - What is the most likely connection between Andrew McCallum and Yiming Yang:



Center-Piece Subgraph (CEPS)

- Given Q query nodes
- Find Center-piece subgraph on b nodes

- Input of **CEPS**:
 - Q Query nodes
 - Budget b
 - k -softAND coefficient



CEPS: 3 steps

■ Individual Score Calculation:

- Measure proximity of each node with respect to individual query node

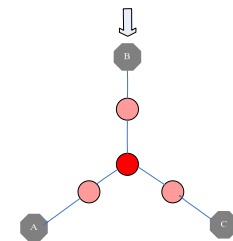
$$r(i, j) = \begin{pmatrix} n \times Q \end{pmatrix}$$

■ Combine Individual Scores:

- Measure importance of a node to the whole query set

$$r(Q, j) = \begin{pmatrix} n \times 1 \end{pmatrix}$$

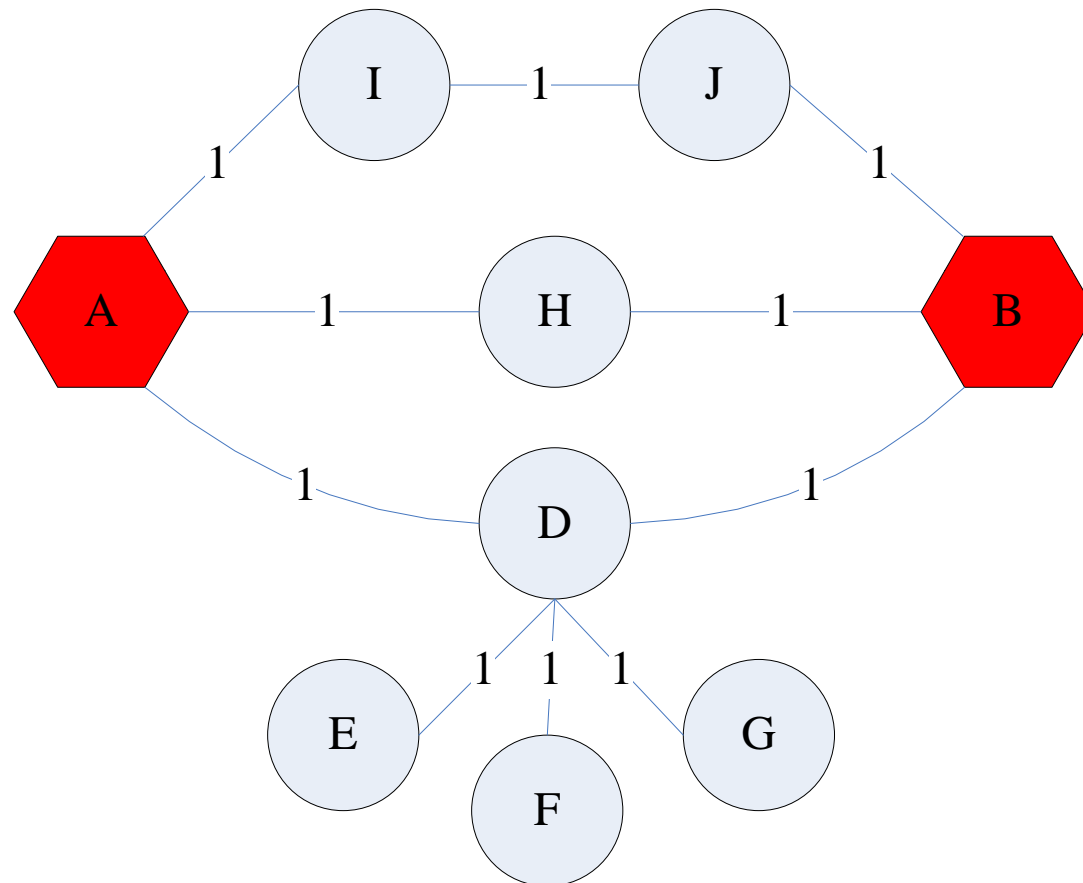
■ “Extract” the connection subgraph: $\arg \max_H g(H)$



Node Score Calculation

- Goal:
 - Calculate importance score $r(i,j) = r_{i,j}$
 - for each **node** j and each **query node** i
- **How to do that?**

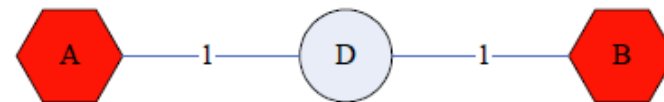
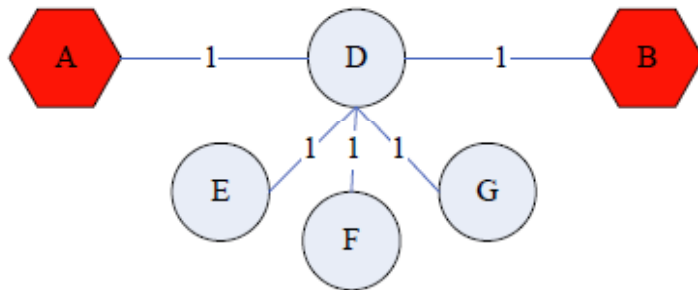
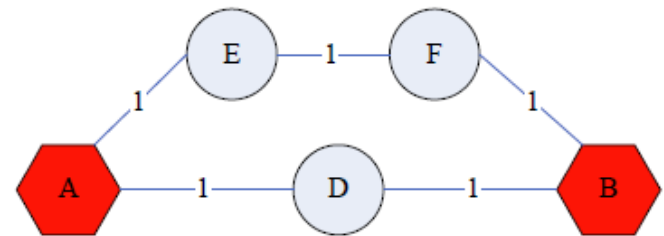
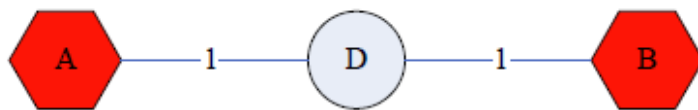
Proximity on Graphs



a.k.a.: Relevance, Closeness, 'Similarity'...

Good proximity measure?

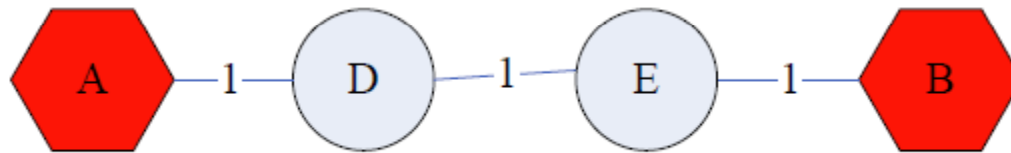
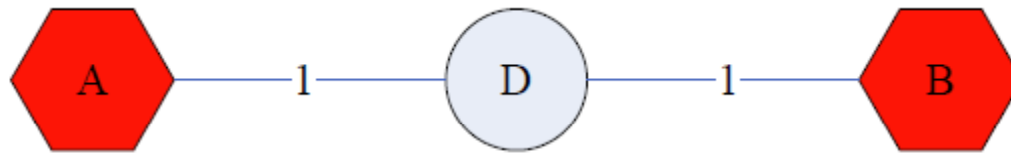
- Shortest path is not good:



- No influence for degree-1 nodes (E, F)!
 - known as 'pizza delivery guy' problem in undirected graph
- Multi-faceted relationships

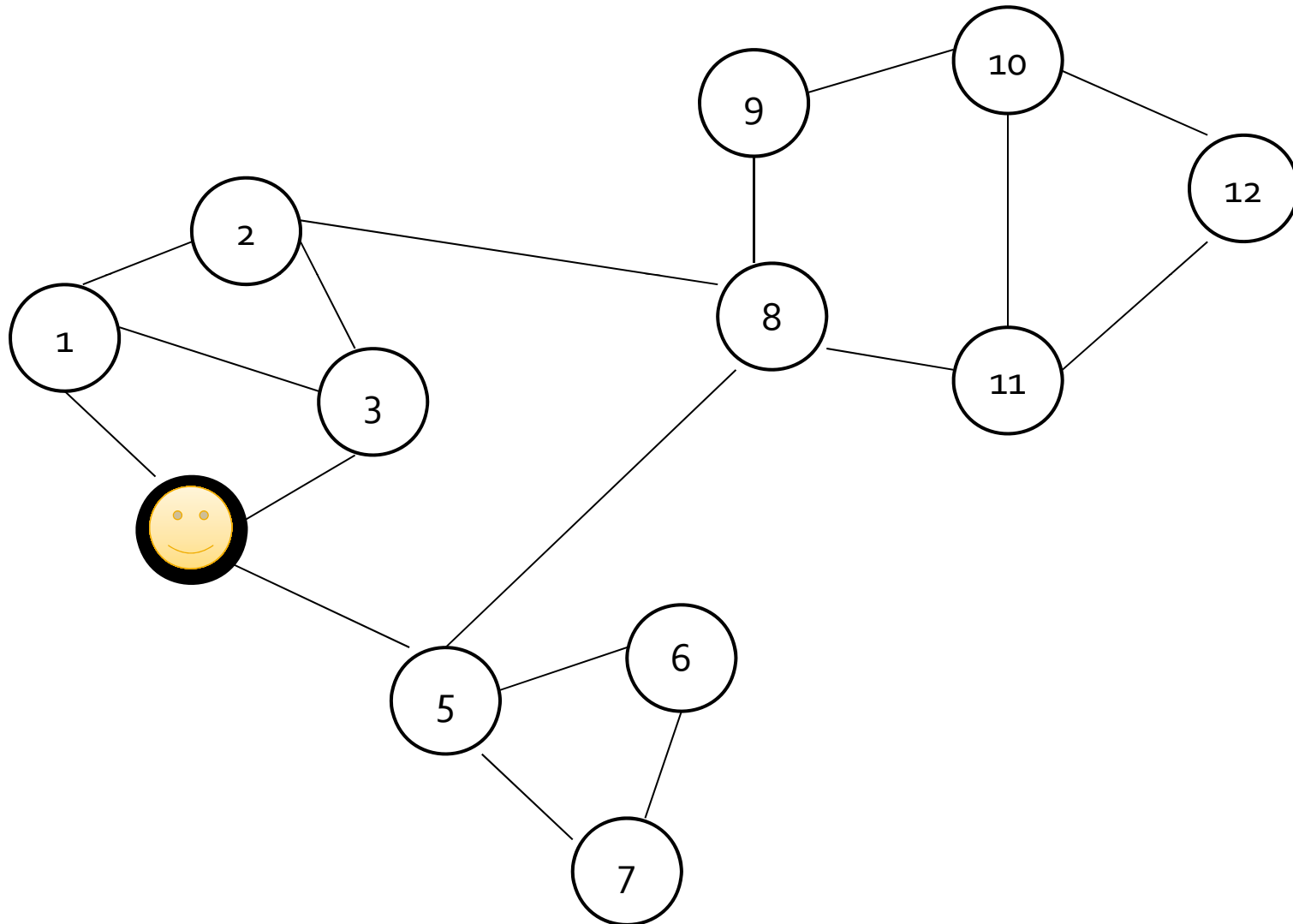
Good proximity measure?

- Max-flow is not good:

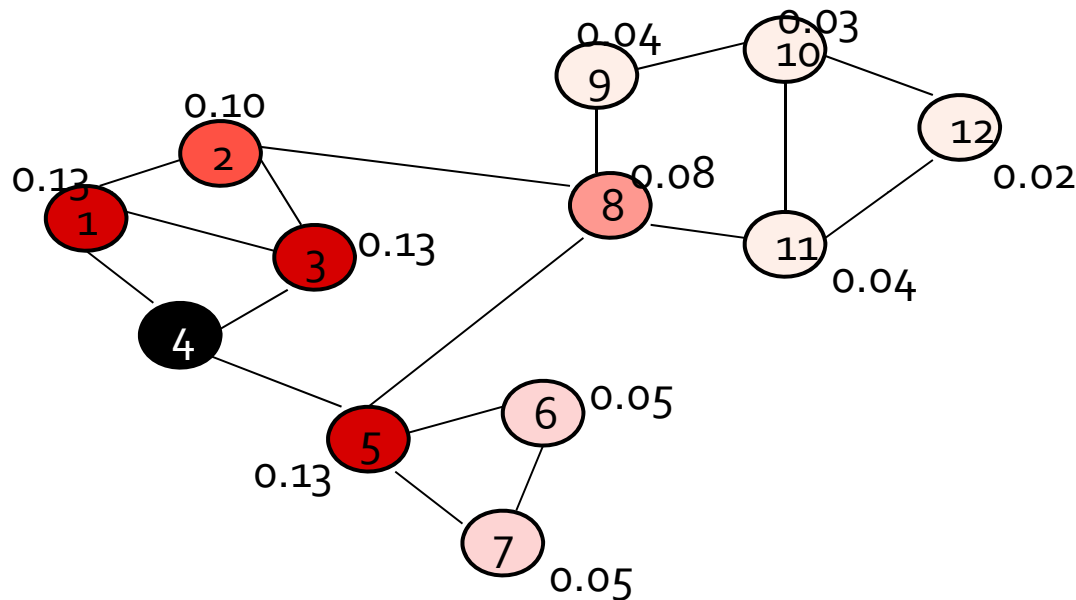


- Does not punish long paths

Random Walk with Restarts



Random Walks with Restarts



Nearby nodes, higher scores

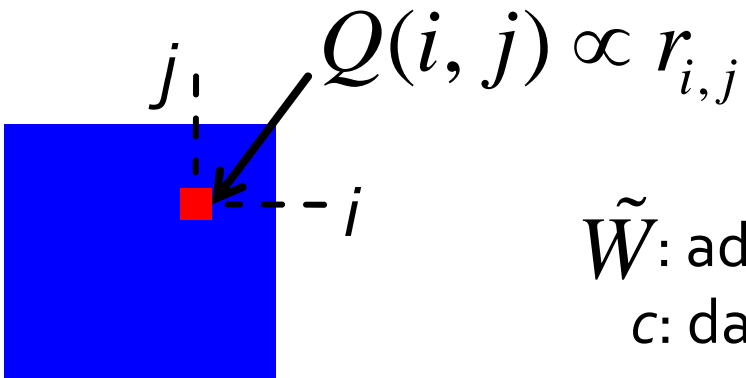
More red, more relevant

	Node 4
Node 1	0.13
Node 2	0.10
Node 3	0.13
Node 4	0.22
Node 5	0.13
Node 6	0.05
Node 7	0.05
Node 8	0.08
Node 9	0.04
Node 10	0.03
Node 11	0.04
Node 12	0.02

Ranking vector

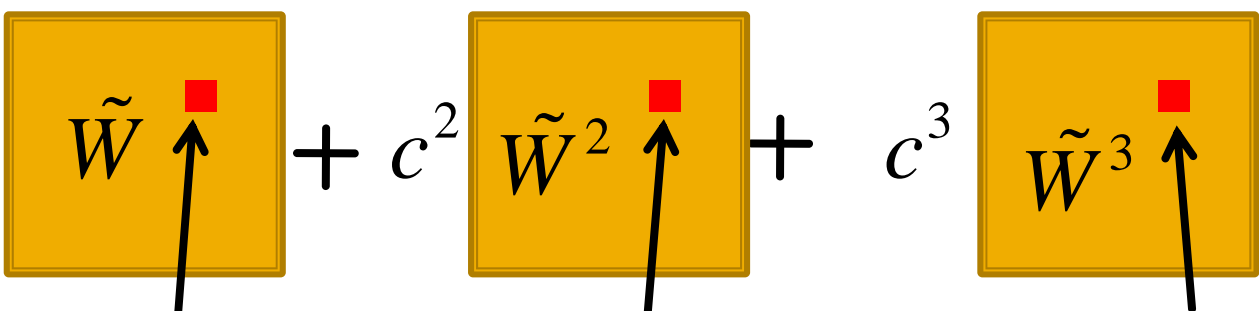
\vec{r}_4

Why is RWR a good score?

$$Q = (I - c\tilde{W})^{-1} =$$


$Q(i, j) \propto r_{i, j}$

\tilde{W} : adjacency matrix.
 c : damping factor

$$Q = I + c \tilde{W} + c^2 \tilde{W}^2 + c^3 \tilde{W}^3 + \dots$$


all paths from i to j with length **1**

all paths from i to j with length **2**

all paths from i to j with length **3**

Computing the score

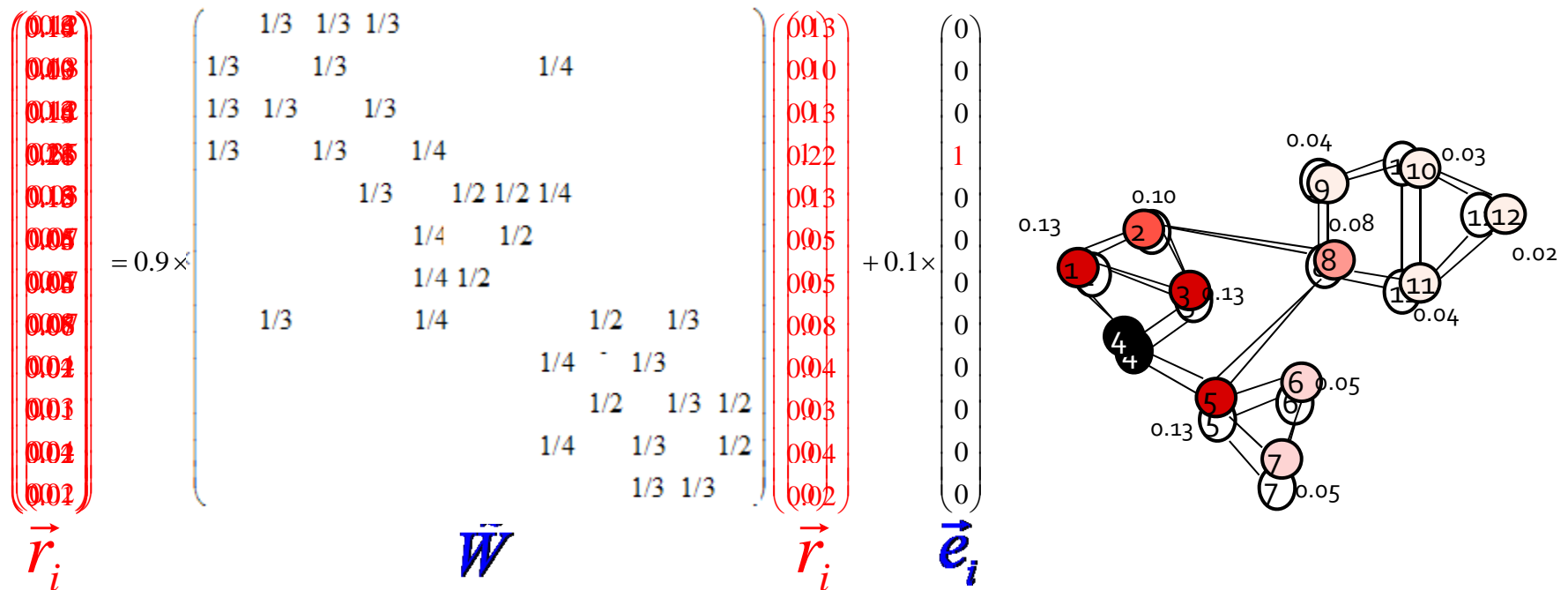
$$\begin{array}{c}
 \left. \begin{array}{c} \left. \begin{array}{c} ? \\ \end{array} \right\} = 0.9 \times \\ \end{array} \right\} \begin{array}{c} \begin{array}{ccccccc} 1/3 & 1/3 & 1/3 & & & & \\ 1/3 & & 1/3 & & & & 1/4 \\ 1/3 & 1/3 & & 1/3 & & & \\ 1/3 & & 1/3 & & 1/4 & & \\ & & & 1/3 & 1/2 & 1/2 & 1/4 \\ & & & 1/4 & & 1/2 & \\ & & & 1/4 & 1/2 & & \\ 1/3 & & & 1/4 & & 1/2 & 1/3 \\ & & & & & 1/4 & 1/3 \\ & & & & & 1/2 & 1/3 & 1/2 \\ 1/4 & & & 1/3 & & 1/2 & & \\ & & & 1/3 & 1/3 & & & \end{array} \\ \end{array} \right\} \begin{array}{c} \left. \begin{array}{c} ? \\ \end{array} \right\} + 0.1 \times \begin{array}{c} \begin{array}{c} 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{array} \\ \end{array}
 \end{array}
 \end{array}$$

Ranking vector
Adjacency matrix
Ranking vector
Starting vector

← Query

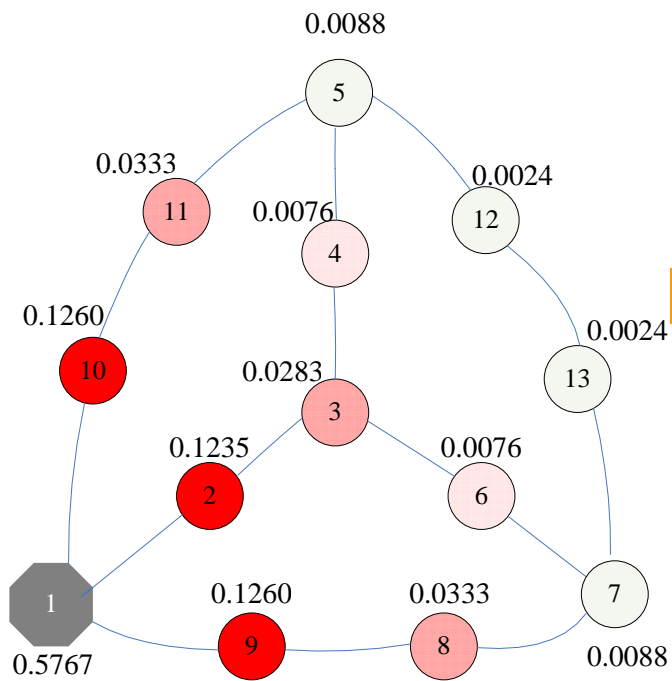
Computing r_i

$$\vec{r}_i[t+1] = c\tilde{W}\vec{r}_i[t] + (1-c)\vec{e}_i$$



No pre-computation / light storage
 Slow on-line response $O(mE)$

Calculating Score $r(i,j)$



	Q1
Node 1	0.5767
Node 2	0.1235
Node 3	0.0283
Node 4	0.0076
Node 5	0.0088
Node 6	0.0076
Node 7	0.0088
Node 8	0.0333
Node 9	0.1260
Node 10	0.1260
Node 11	0.0333
Node 12	0.0024
Node 13	0.0024

Combining the scores

- The RWR score $r(i,j)$
 - i ... query node, $i \in Q$
 - j ... node in the network
- Combine $r(i,j)$ to get the score $r(Q,j)$
 - **1) AND query:** Meeting probability
 - prob. that $|Q|$ random particles coincide on node j

$$r(Q, j) \triangleq r(Q, j, Q) = \prod_{i=1}^Q r(i, j)$$

Combining the scores: k-SoftAND

- OR query:

$$r(Q, j) \triangleq r(Q, j, Q) = \prod_{i=1}^Q r(i, j)$$

- At least one particle arrives to node j :
(i.e., node j is important to at least 1 query node)

$$r(Q, j) \triangleq r(Q, j, 1) = 1 - \prod_{i=1}^Q (1 - r(i, j))$$

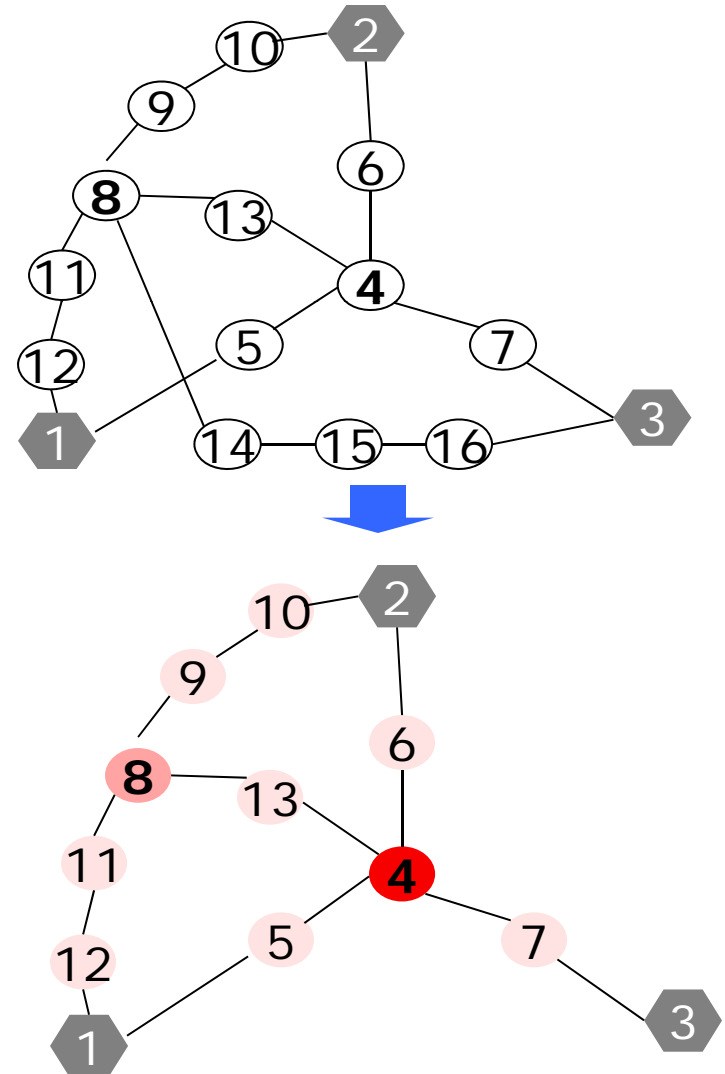
- k-SoftAND query:

- Node j is important to at least k query nodes
- Can be computed recursively:

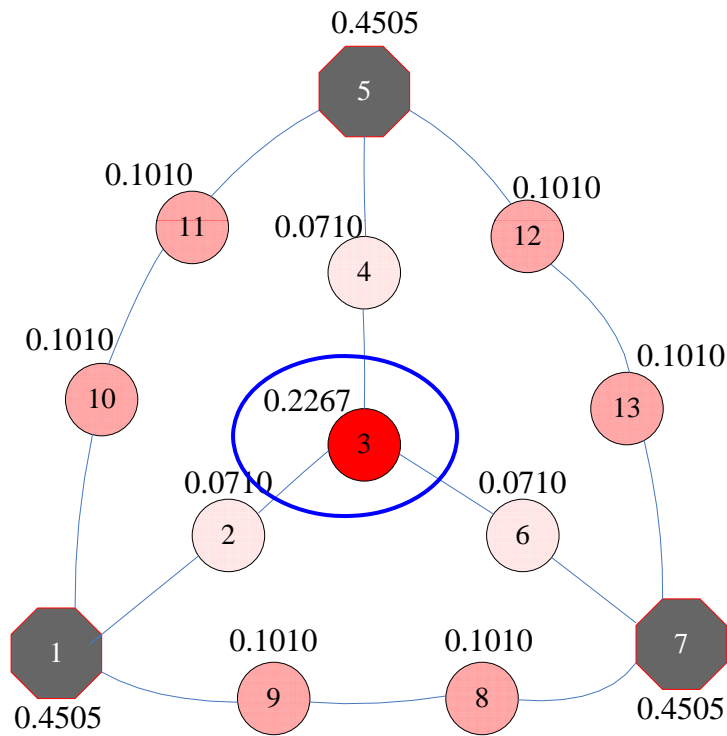
$$r(Q, j, k) = r(Q', j, k - 1) \cdot r(Q, j) + r(Q', j, k) \cdot (1 - r(Q, j))$$

Extracting Connection Subgraph

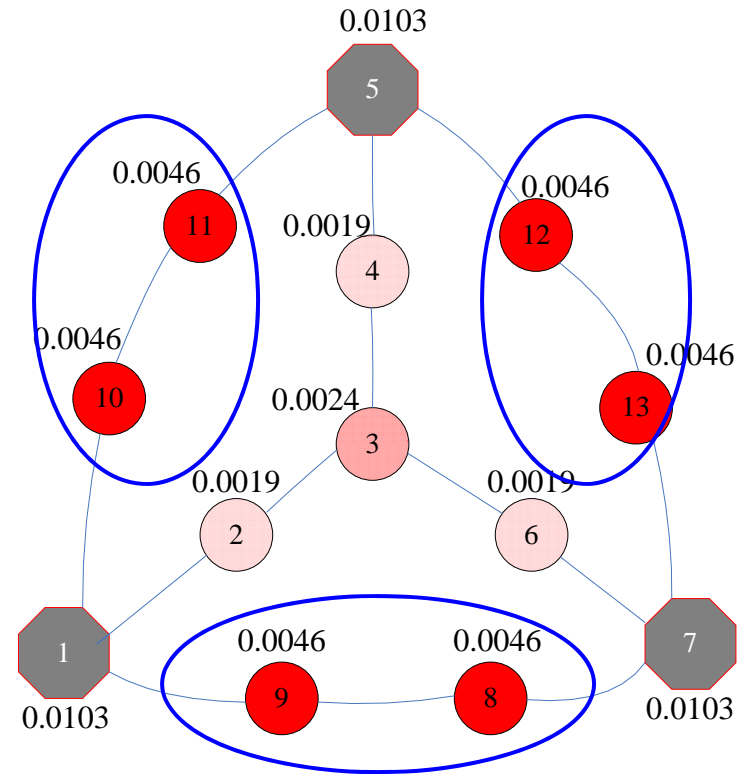
- Goal:
 - Extract a sub-graph S on b nodes that maximizes the $\sum_{u \in S} r(Q, u)$
- Idea:
 - Iterate until budget is reached
 - Pick not-yet-selected node $j = \arg \max_j r(Q, j)$
 - Find good paths to all query nodes
 - Add the paths to S



AND vs. k-SoftAND

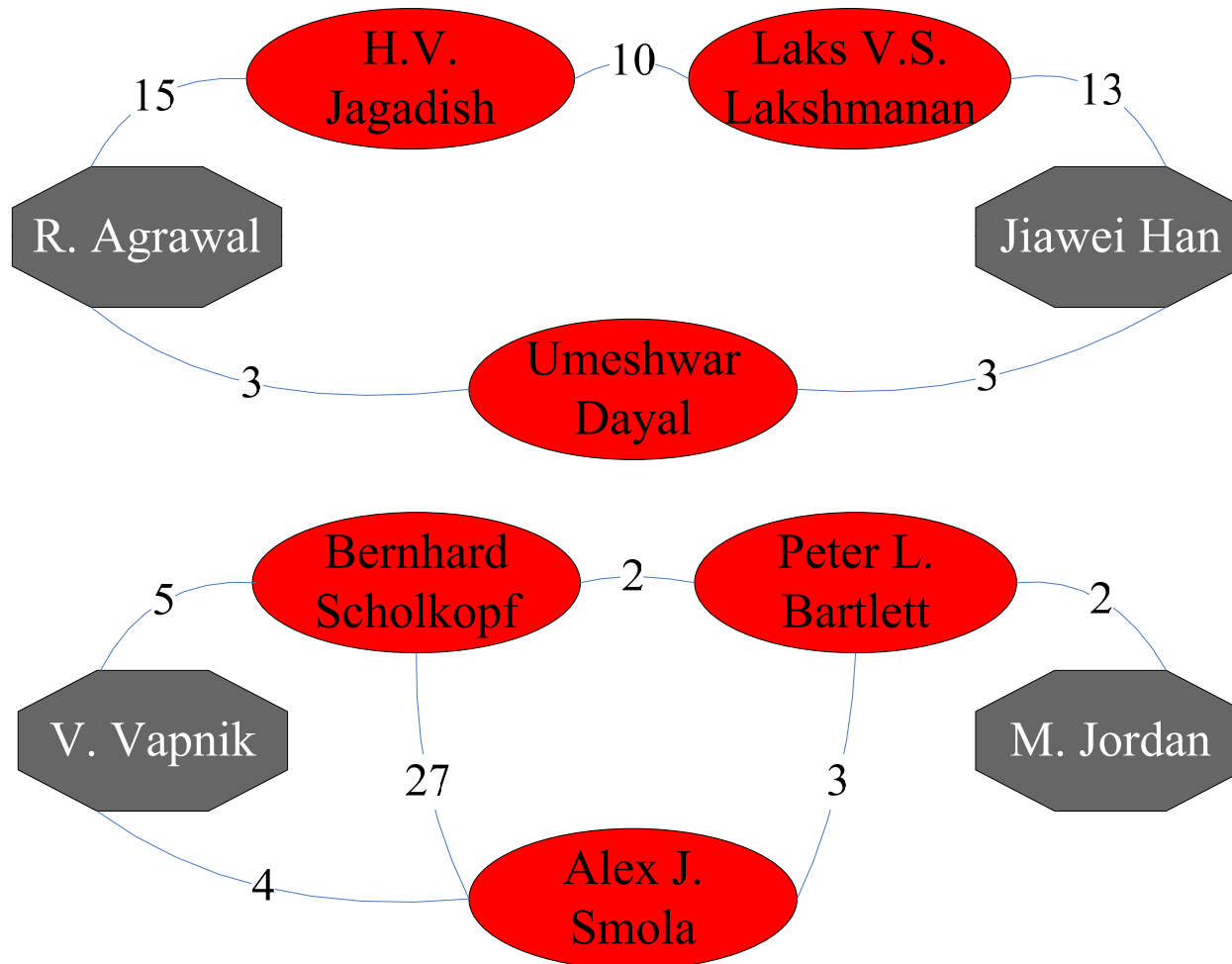


AND



2-SoftAND

Example:



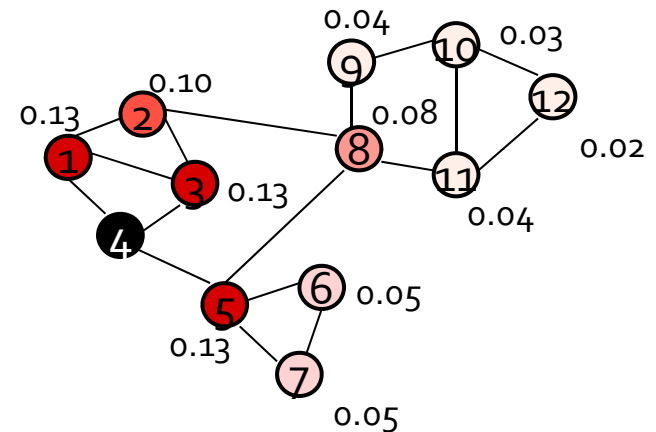
Speeding up the response time

- Solving the random walk for each query separately is time consuming
 - Not appropriate for real-time
- Can we do better?
- Idea 1) Pre-compute scores for each possible query node i

Pre-compute all the scores

Q^{-1} .

r_1	r_2	r_3	r_4	r_5	r_6	r_7	r_8	r_9	r_{10}	r_{11}	r_{12}
0.20	0.13	0.14	0.13	0.68	0.56	0.56	0.63	0.44	0.35	0.39	0.34
0.28	0.20	0.13	0.96	0.64	0.53	0.53	0.85	0.60	0.48	0.53	0.45
0.14	0.13	0.20	1.29	0.68	0.56	0.56	0.63	0.44	0.35	0.39	0.33
0.13	0.10	0.13	2.06	0.95	0.78	0.78	0.61	0.43	0.34	0.38	0.32
0.09	0.09	0.09	1.27	2.41	1.97	1.97	1.05	0.73	0.58	0.66	0.56
0.03	0.04	0.04	0.52	0.98	2.06	1.37	0.43	0.30	0.24	0.27	0.22
0.03	0.04	0.04	0.52	0.98	1.37	2.06	0.43	0.30	0.24	0.27	0.22
0.08	0.11	0.04	0.82	1.05	0.86	0.86	2.13	1.49	1.19	1.33	1.13
0.03	0.04	0.03	0.28	0.36	0.30	0.30	0.74	1.78	1.00	0.76	0.79
0.04	0.04	0.04	0.34	0.44	0.36	0.36	0.89	1.50	2.45	1.54	1.80
0.04	0.05	0.04	0.38	0.49	0.40	0.40	1.00	1.14	1.54	2.28	1.72
0.02	0.03	0.02	0.21	0.28	0.22	0.22	0.56	0.79	1.20	1.14	2.05



$$Q^{-1} = (I - c\tilde{W})^{-1}$$

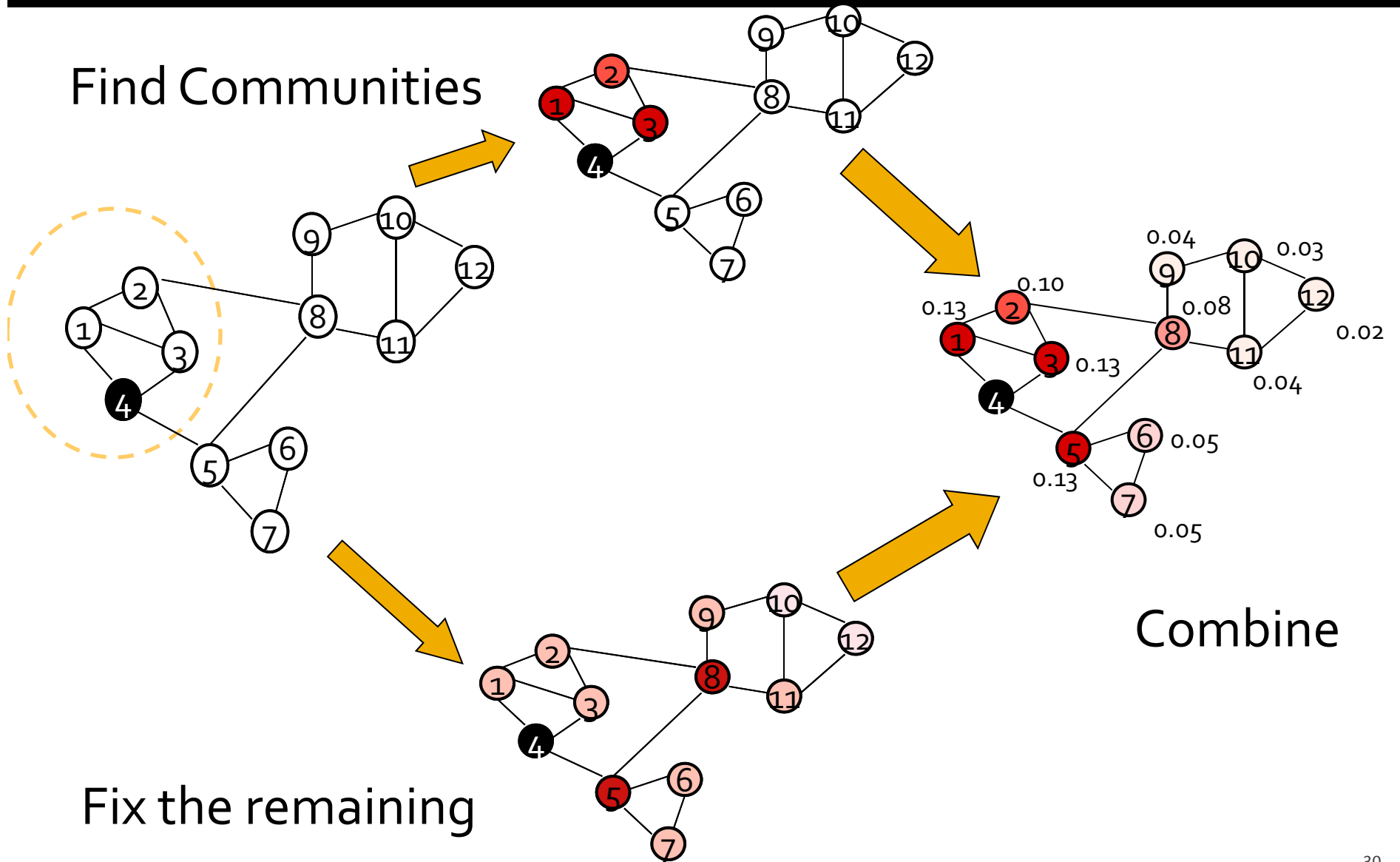
Fast on-line response

Heavy pre-computation/storage cost

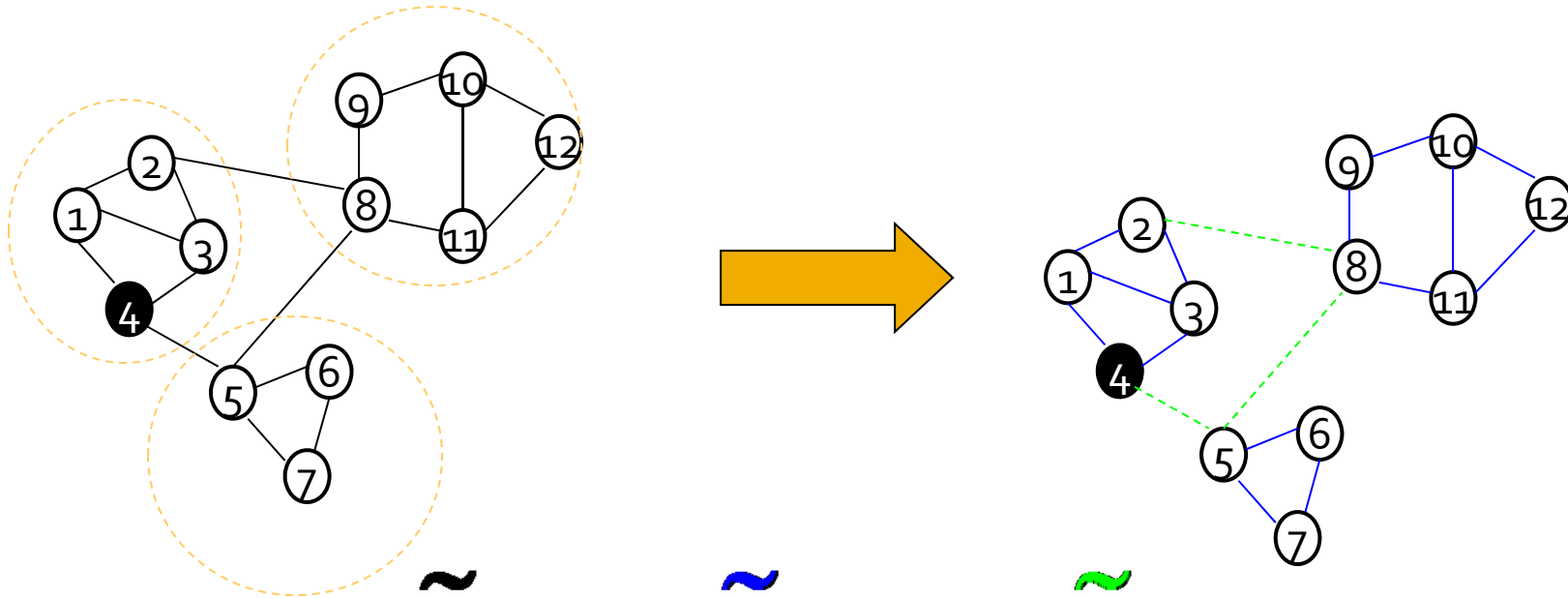
$O(n^3)$

$O(n^2)$

Idea 2: Approximation



Step 1: Partition the graph

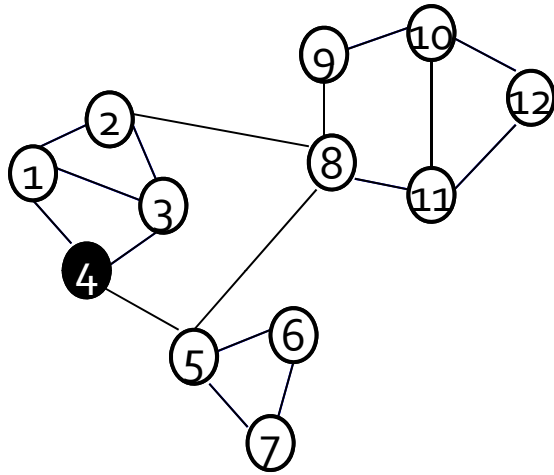


$$\tilde{W} = \tilde{W}_1 + \tilde{W}_2$$

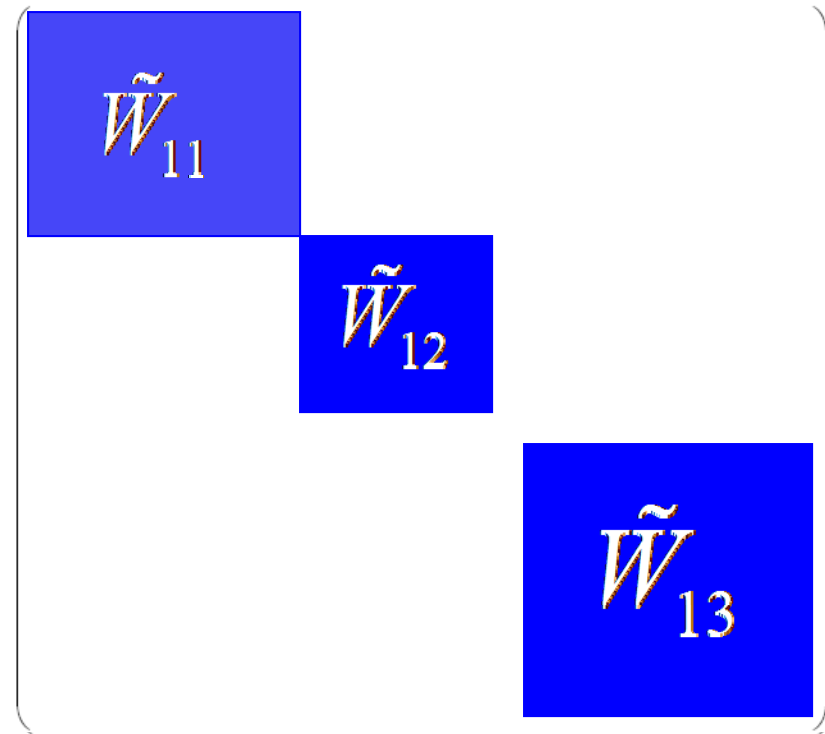
Within-partition links

cross-partition links

Step 2: Make W_1 block-diagonal

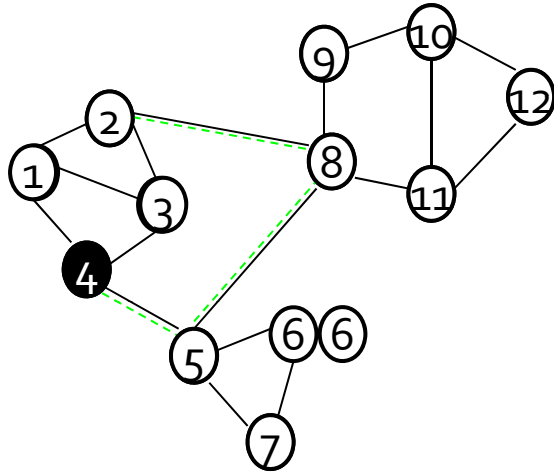


$$\tilde{W}_1 = \begin{pmatrix} \tilde{W}_{11} & & \\ & \tilde{W}_{22} & \\ & & \tilde{W}_{33} \end{pmatrix} =$$



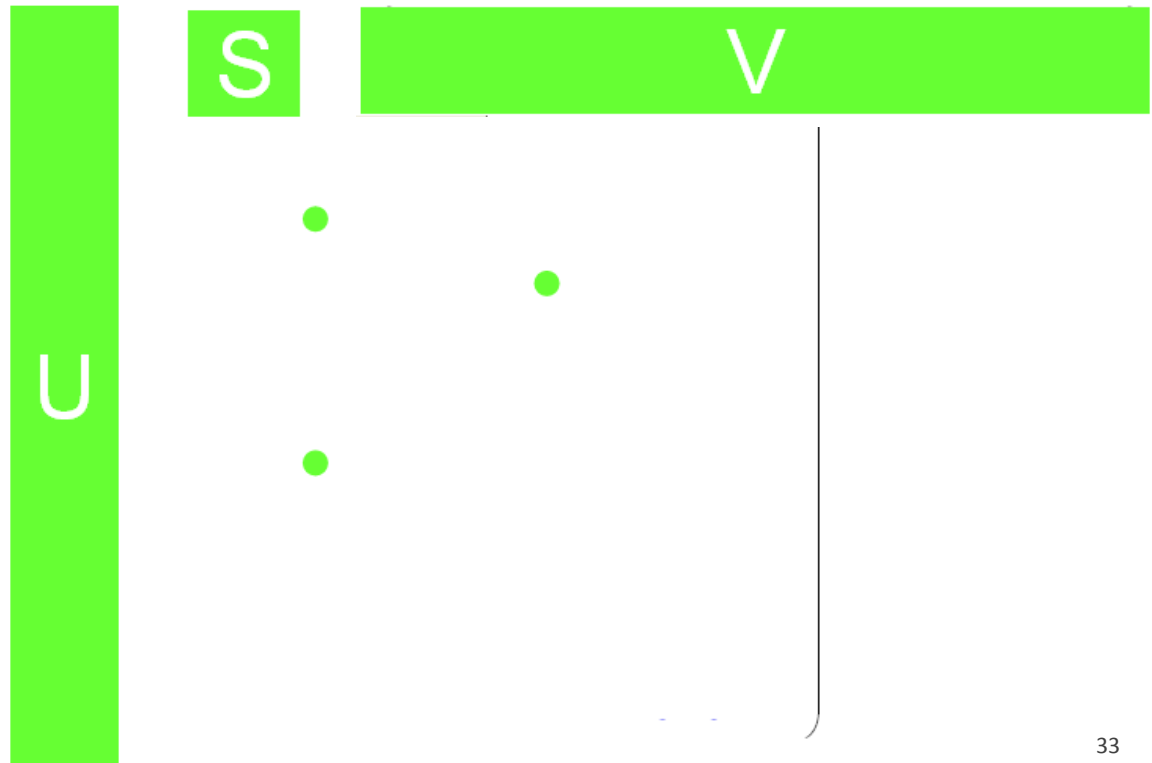
$$Q_{1,i}^{-1} = (\mathbf{I} - \tilde{W}_{1,i})^{-1}$$

Step 3: Cross community links



$$|S| \ll |\tilde{W}_2|$$

$$\tilde{W}_2 \approx USV =$$



Final algorithm: Offline strage

- Offline stage:

$$Q_1^{-1} = \begin{pmatrix} Q_{1,1} & & & \\ & Q_{1,2} & & \\ & & \dots & \\ & & & Q_{1,k} \end{pmatrix}$$

Communities

$$\tilde{\Lambda} = \begin{pmatrix} 0.39 & -0.31 & 0.03 & -0.03 \\ -0.28 & 1.18 & -0.10 & 0.06 \\ -0.00 & -0.05 & 0.17 & 0.03 \\ -0.03 & 0.01 & 0.03 & 0.19 \end{pmatrix}$$

$$\tilde{\Lambda} = (S^{-1} - cVQ_1^{-1}U)^{-1}$$

Bridges

$$\tilde{W} = \tilde{W}_1 + \tilde{W}_2$$

One more thing:

- Sherman-Morrison-Woodbury matrix identity:
 - Inverse of a rank- k correction of matrix X can be computed by doing a rank- k correction to X^{-1} :

$$(\mathbf{X} - \mathbf{U}\mathbf{S}\mathbf{V})^{-1} = \mathbf{X}^{-1} + \mathbf{X}^{-1}\mathbf{U}\tilde{\mathbf{\Lambda}}\mathbf{V}\mathbf{X}^{-1}$$

$$\text{where } \tilde{\mathbf{\Lambda}} = (\mathbf{S}^{-1} - \mathbf{V}\mathbf{X}^{-1}\mathbf{U})^{-1}$$

- In our case:

$$\mathbf{Q}^{-1} = c(\mathbf{I} - \mathbf{W})^{-1} = c(\mathbf{I} - \mathbf{W}_1 - \mathbf{W}_2)^{-1} =$$

$$\mathbf{Q}^{-1} \approx \mathbf{Q}_1^{-1} + c\mathbf{Q}_1^{-1}\mathbf{U}\tilde{\mathbf{\Lambda}}\mathbf{V}\mathbf{Q}_1^{-1}$$

Final algorithm: Online stage

- Online stage:
 - Compute column i of:

$$Q^{-1} \approx Q_1^{-1} + cQ_1^{-1}U\tilde{\Lambda}VQ_1^{-1}$$

- Using:

$$\vec{r}_i = (1 - c)(Q_1^{-1}e_i + cQ_1^{-1}U\tilde{\Lambda}VQ_1^{-1}e_i)$$

Acknowledgment

- Most slides are borrowed from the tutorial by Hanghang Tong and Christos Faloutsos from CMU
- http://www.cs.cmu.edu/~htong/tut/cikm2008/cikm_tutorial.html