

# Network Science

## #6 Insights on PageRank

© 2020 T. Erseghe



# Measuring closeness with PageRank

# How can we use PageRank?

Want to know about a specific topic? **TopicSpecific** PageRank

Want to measure proximity/similarity to a node? **Local PageRank**

... appropriately select your teleport vector  $q$  !



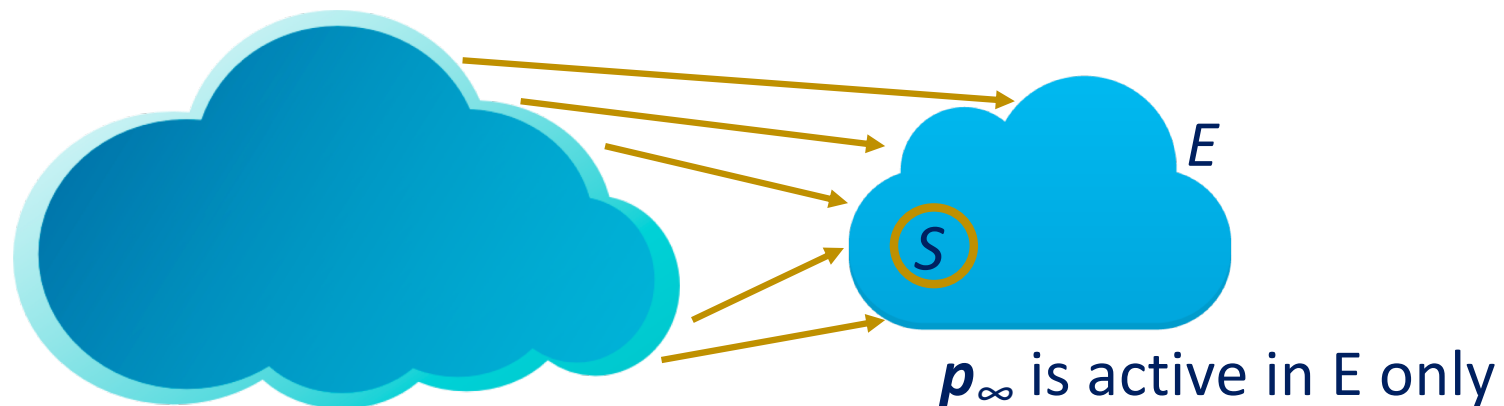
# Topic specific PageRank

## Idea

- ❑ Bias the random walk towards a **topic specific teleport set**  $S$  of nodes, i.e., make sure that  $q$  is active in  $S$  only
- ❑  $S$  should contain only pages that are relevant to the topic

## Result

- ❑ The random walk **deterministically** ends in a small set  $E$ , containing  $S$ , and being in some sense close to it



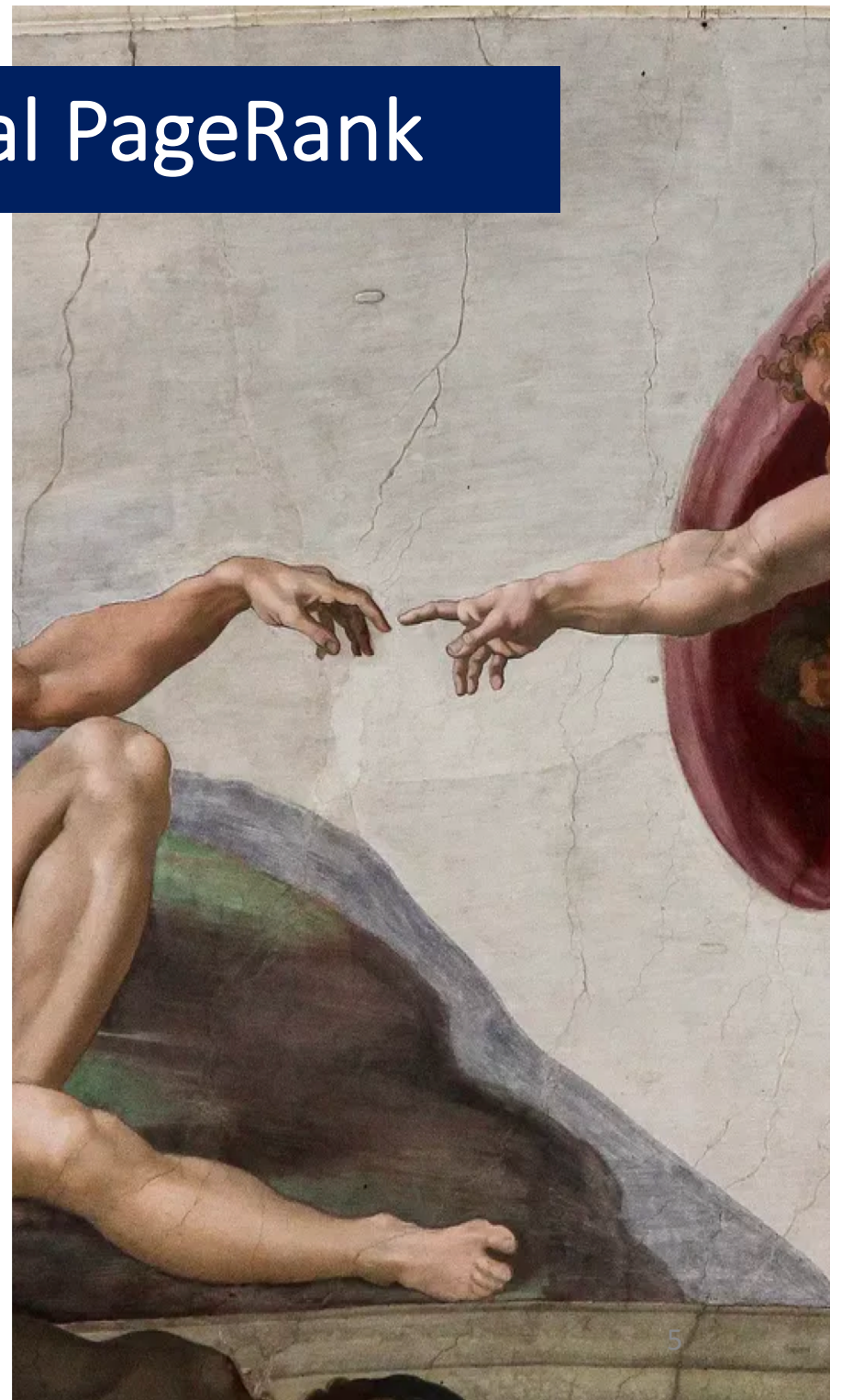
# Measuring closeness: Local PageRank

## Idea

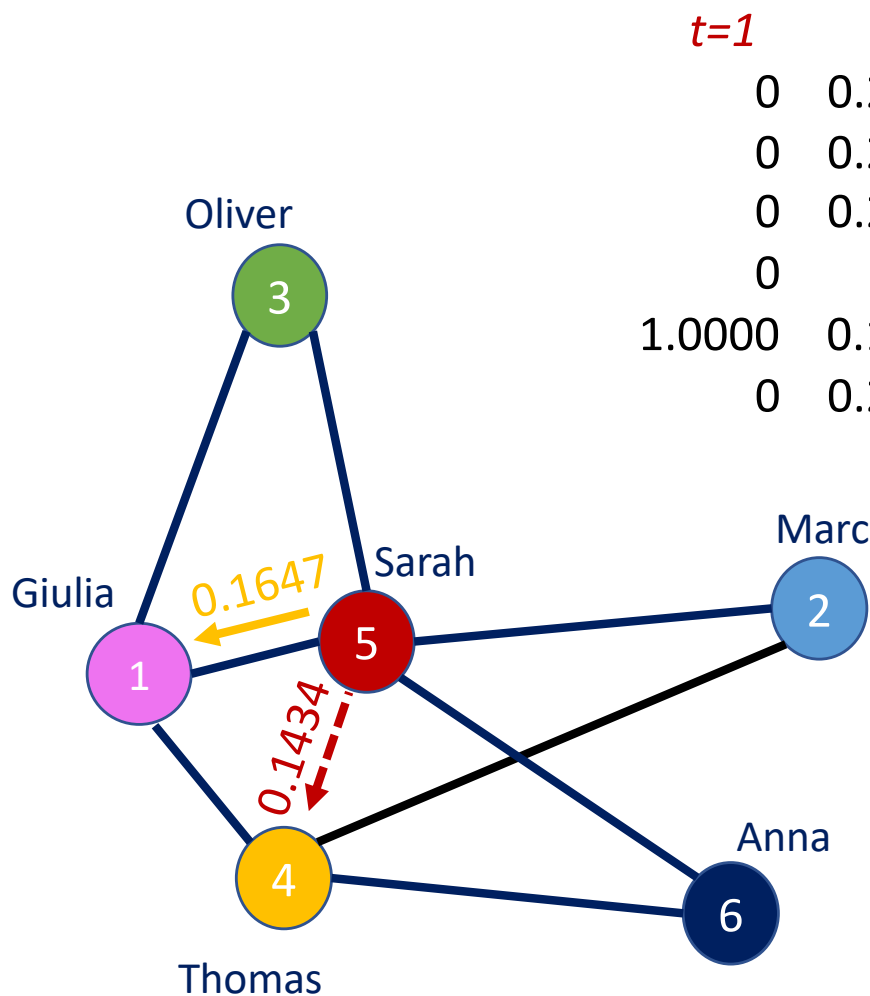
- ❑ Measure **similarity** / closeness to node  $i$  by applying TopicSpecific PageRank with teleport set  $S=\{i\}$

## Result

- ❑ Measures direct and indirect multiple connections, their quality, degree or weight



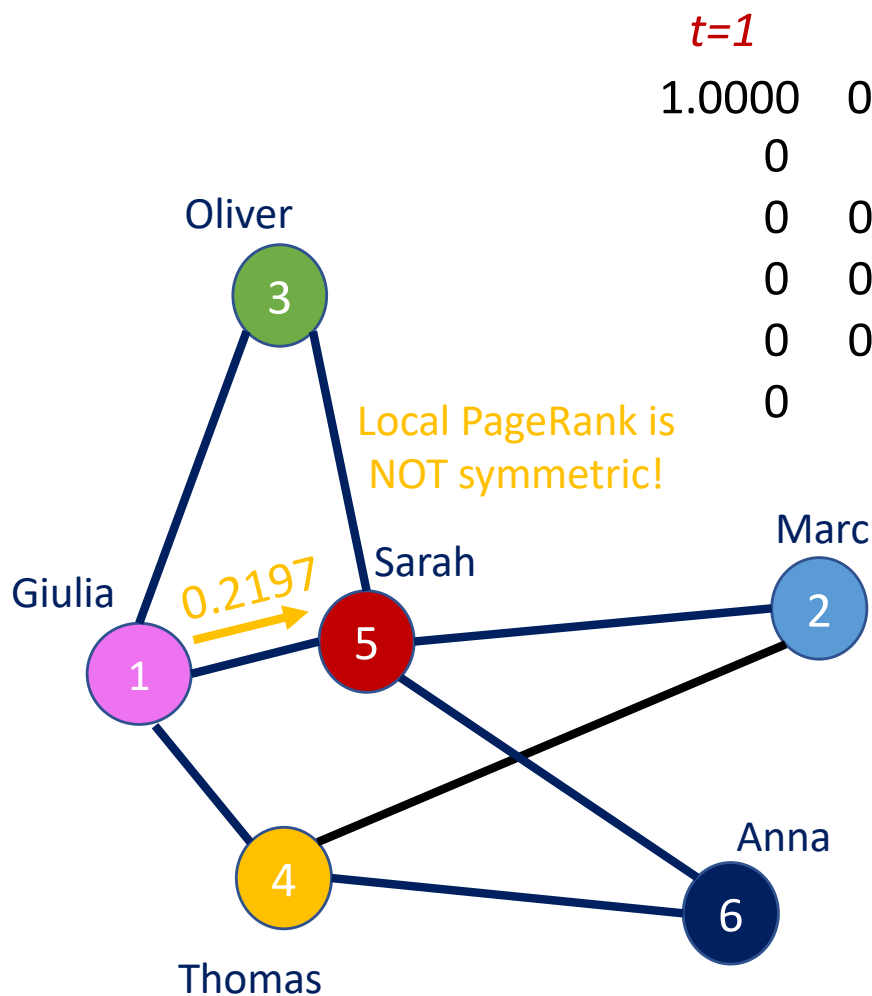
# Example: who's Sarah's best friend?



|        | <i>t=1</i> | <i>2</i> | <i>3</i> | <i>4</i> | <i>5</i> |
|--------|------------|----------|----------|----------|----------|
| Oliver | 0          | 0.2125   | 0.1222   | 0.2096   | 0.1290   |
| Giulia | 0          | 0.2125   | 0.0319   | 0.1705   | 0.0708   |
| Sarah  | 0          | 0.2125   | 0.0921   | 0.1369   | 0.1127   |
| Marc   | 0          | 0        | 0.2408   | 0.0617   | 0.2043   |
| Thomas | 1.0000     | 0.1500   | 0.4811   | 0.2508   | 0.4125   |
| Anna   | 0          | 0.2125   | 0.0319   | 0.1705   | 0.0708   |

|        | <i>10</i> | <i>20</i> | <i>50</i> | <i>75</i> | <i>100</i> |        |
|--------|-----------|-----------|-----------|-----------|------------|--------|
| 0.1743 | 0.1653    | 0.1647    | 0.1647    | 0.1647    | 0.1647     | Giulia |
| 0.1238 | 0.1144    | 0.1138    | 0.1138    | 0.1138    | 0.1138     | Marc   |
| 0.1206 | 0.1199    | 0.1199    | 0.1199    | 0.1199    | 0.1199     | Oliver |
| 0.1285 | 0.1426    | 0.1434    | 0.1434    | 0.1434    | 0.1434     | Thomas |
| 0.3290 | 0.3435    | 0.3444    | 0.3444    | 0.3444    | 0.3444     | Sarah  |
| 0.1238 | 0.1144    | 0.1138    | 0.1138    | 0.1138    | 0.1138     | Anna   |

# Example: who's Giulia's best friend?

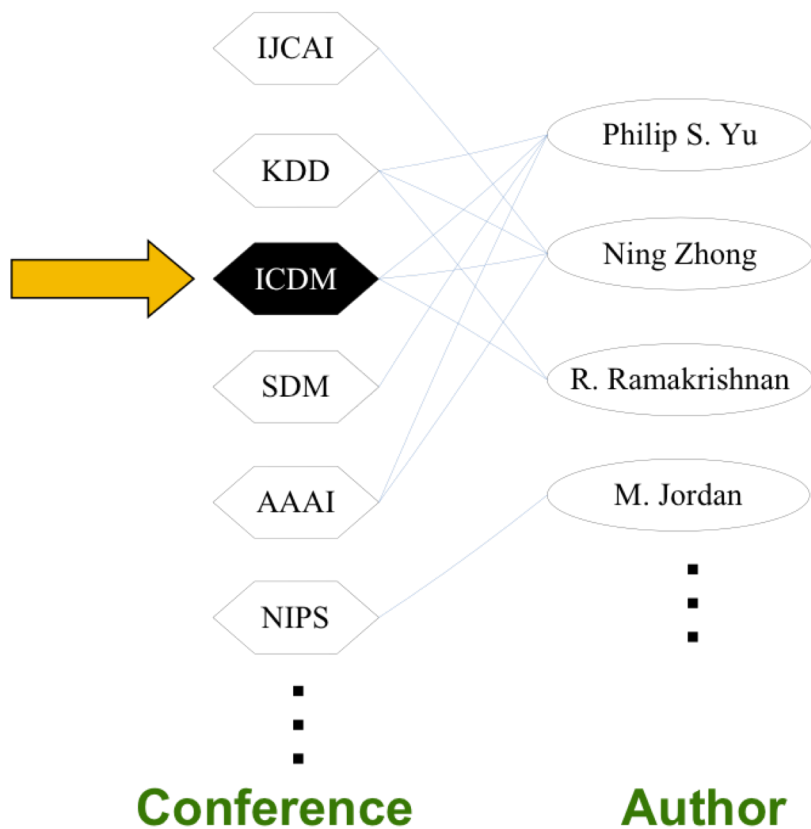


|        | <i>t=1</i> | <i>2</i> | <i>3</i> | <i>4</i> | <i>5</i> |
|--------|------------|----------|----------|----------|----------|
| Giulia | 1.0000     | 0.1500   | 0.4109   | 0.2403   | 0.3404   |
| Oliver | 0          | 0        | 0.1405   | 0.0467   | 0.1262   |
| Sarah  | 0          | 0.2833   | 0.1027   | 0.1510   | 0.1275   |
| Marc   | 0          | 0.2833   | 0.0425   | 0.2358   | 0.1078   |
| Thomas | 0          | 0.2833   | 0.1629   | 0.2795   | 0.1719   |
| Anna   | 0          | 0        | 0.1405   | 0.0467   | 0.1262   |

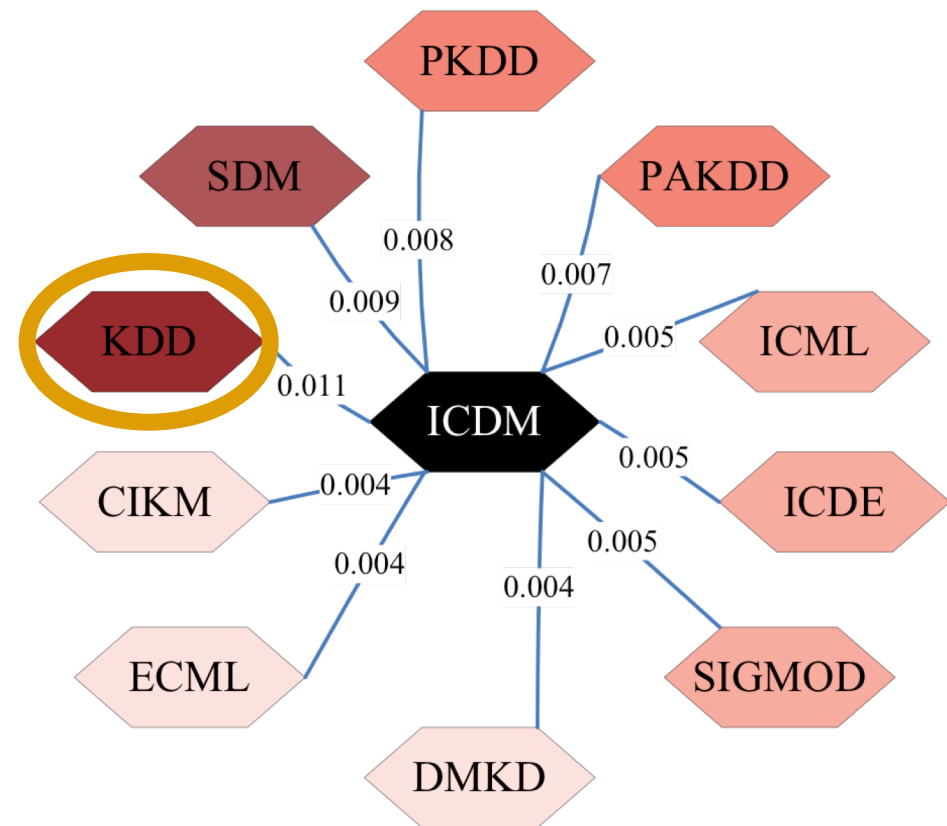
|        | <i>10</i> | <i>20</i> | <i>50</i> | <i>75</i> | <i>100</i> |        |
|--------|-----------|-----------|-----------|-----------|------------|--------|
| Giulia | 0.2909    | 0.2985    | 0.2989    | 0.2989    | 0.2989     | Giulia |
| Marc   | 0.0848    | 0.0926    | 0.0931    | 0.0931    | 0.0931     | Marc   |
| Oliver | 0.1309    | 0.1313    | 0.1314    | 0.1314    | 0.1314     | Oliver |
| Thomas | 0.1763    | 0.1645    | 0.1638    | 0.1638    | 0.1638     | Thomas |
| Sarah  | 0.2324    | 0.2204    | 0.2197    | 0.2197    | 0.2197     | Sarah  |
| Anna   | 0.0848    | 0.0926    | 0.0931    | 0.0931    | 0.0931     | Anna   |

# Example

What is the most related conference to ICDM?



Top 10 ranking results

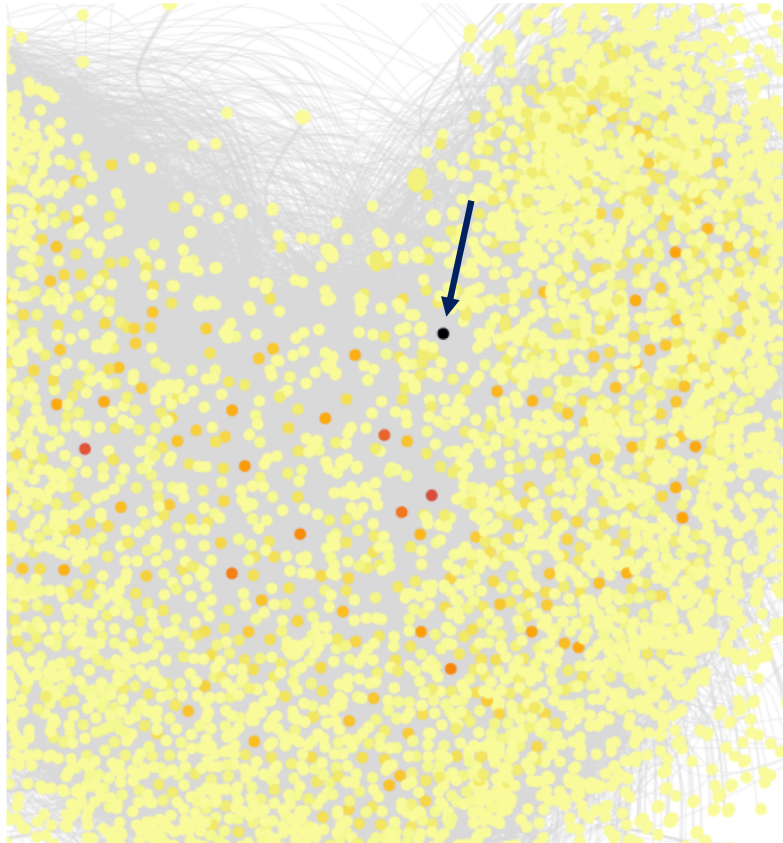


ICDM = international conf. on data mining  
KDD = knowledge discovery and data mining



# Local PageRank (authorities , A)

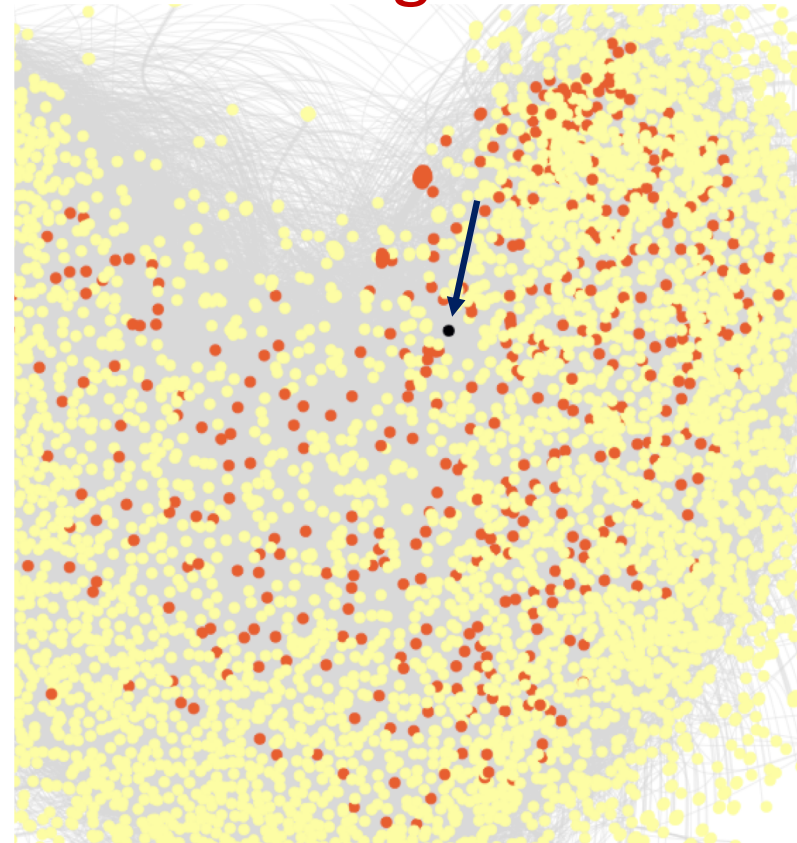
Local PageRank



neighbours authority score =  
local node  $\rightarrow$  neighbours

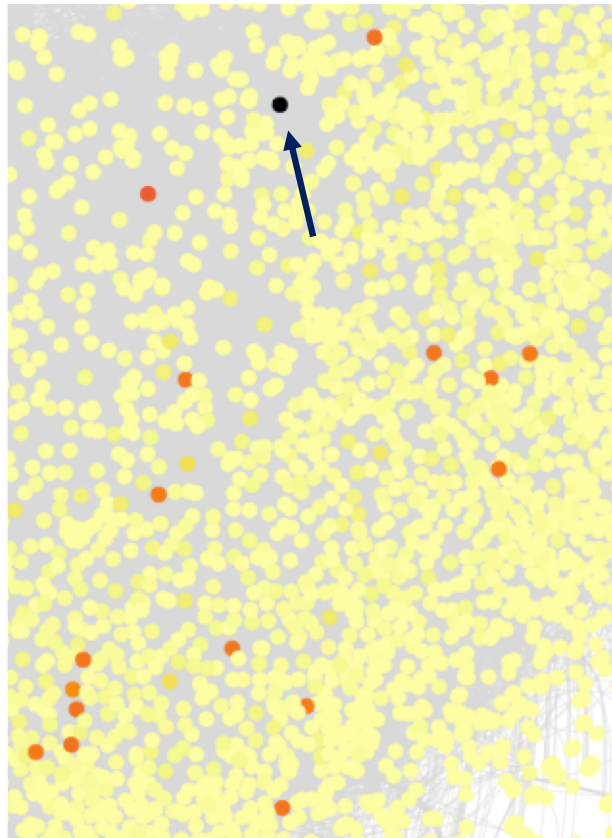
1-hop

out-neighbours



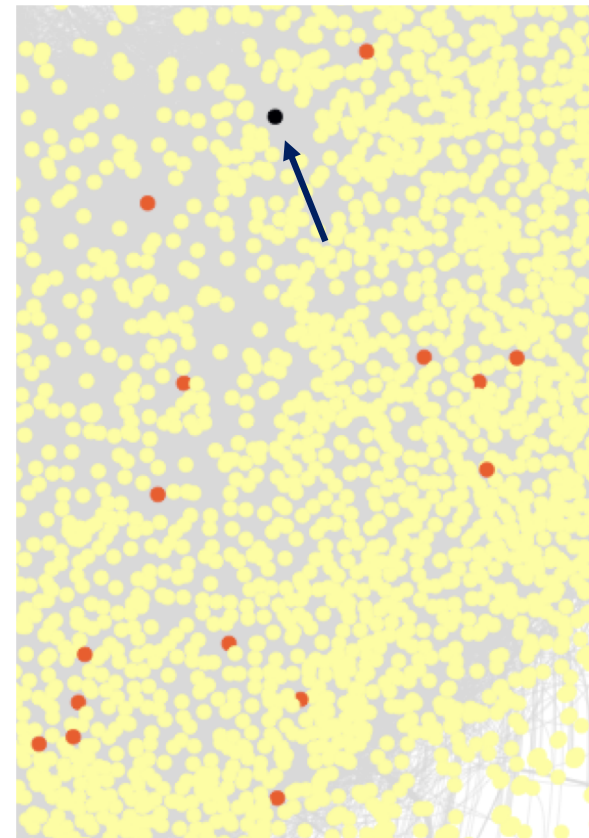
# Local PageRank (hubs, $A^T$ )

Local PageRank



neighbours **hub score** =  
neighbours  $\rightarrow$  local node

1-hop  
**in-neighbours**



# Combating spam farms

# Spam farms



Google™ [Web](#) [Images](#) [Groups](#) [News](#) [Froogle](#) [Local](#) [more »](#)

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

**Web** Results 1 - 10 of about 969,000 for [miserable failure](#). (0.06 seconds)

[Biography of President George W. Bush](#)  
Biography of the president from the official White House web site.  
[www.whitehouse.gov/president/gwbbio.html](http://www.whitehouse.gov/president/gwbbio.html) - 29k - [Cached](#) - [Similar pages](#)  
[Past Presidents](#) - [Kids Only](#) - [Current News](#) - [President](#)  
[More results from www.whitehouse.gov »](#)

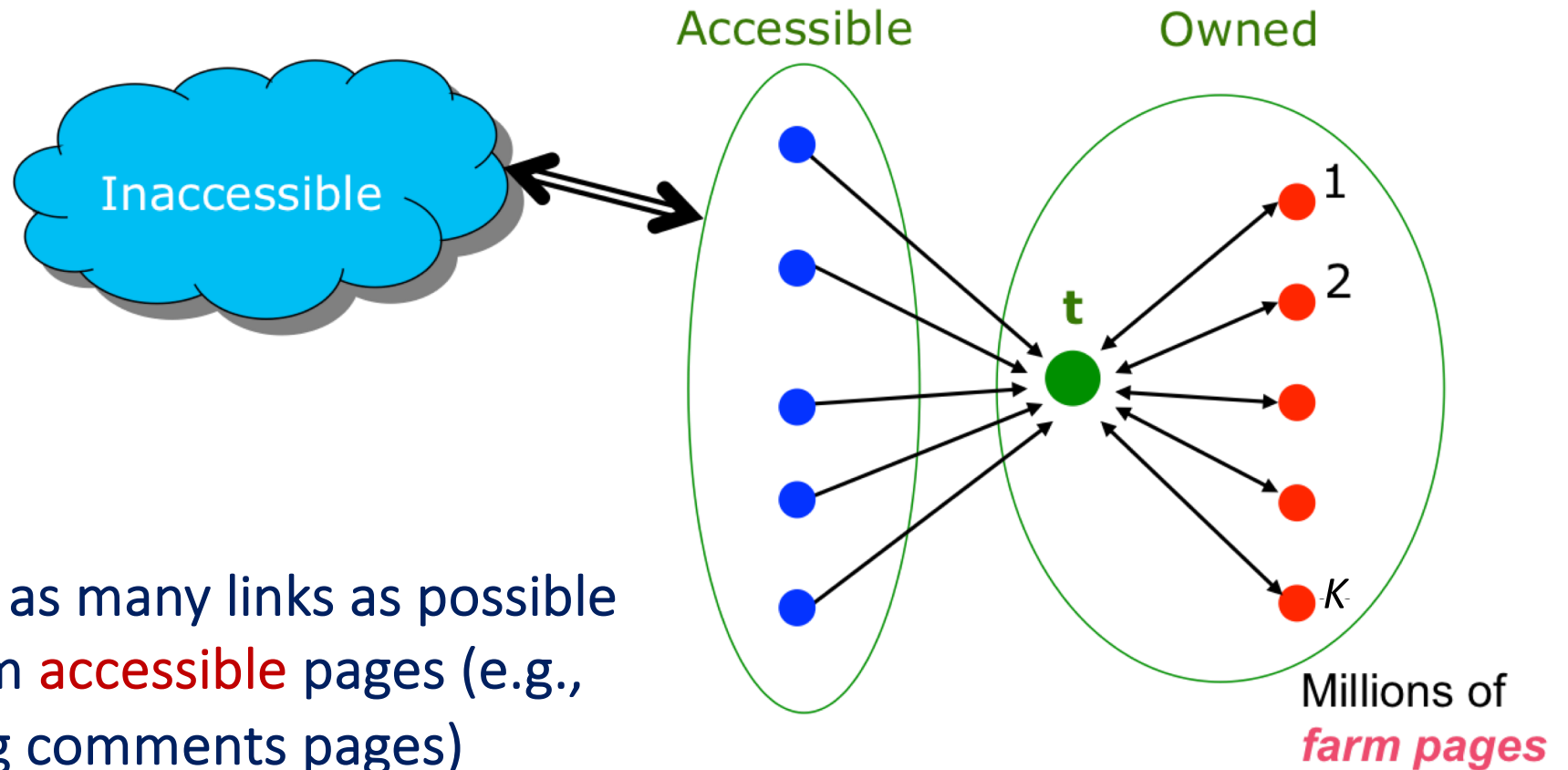
[Welcome to MichaelMoore.com!](#)  
Official site of the gadfly of corporations, creator of the film Roger and Me and the television show The Awful Truth. Includes mailing list, message board, ...  
[www.michaelmoore.com/](http://www.michaelmoore.com/) - 35k - Sep 1, 2005 - [Cached](#) - [Similar pages](#)

[BBC NEWS | Americas | 'Miserable failure' links to Bush](#)  
Web users manipulate a popular search engine so an unflattering description leads to the president's page.  
[news.bbc.co.uk/2/hi/americas/3298443.stm](http://news.bbc.co.uk/2/hi/americas/3298443.stm) - 31k - [Cached](#) - [Similar pages](#)

[Google's \(and Inktomi's\) Miserable Failure](#)  
A search for **miserable failure** on Google brings up the official George W. Bush biography from the US White House web site. Dismissed by Google as not a ...  
[searchenginewatch.com/sereport/article.php/3296101](http://searchenginewatch.com/sereport/article.php/3296101) - 45k - Sep 1, 2005 - [Cached](#) - [Similar pages](#)

Google bombs in the  
2004 elections

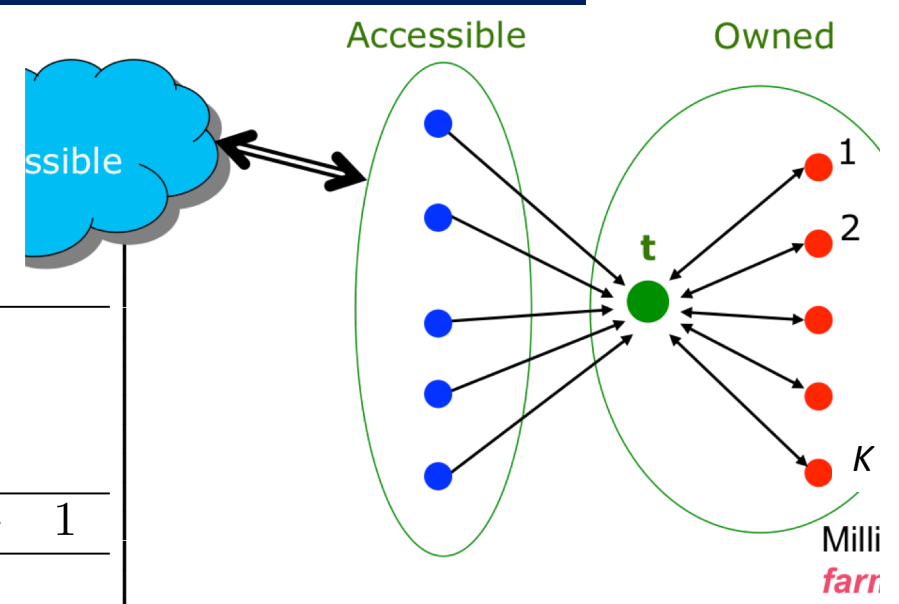
# Spamming method



1. Get as many links as possible from **accessible** pages (e.g., blog comments pages)
2. Construct **link farm** to get a PageRank multiplier effect

# PageRank analysis

|        |                   |                   |                                     |                   |
|--------|-------------------|-------------------|-------------------------------------|-------------------|
| inacc. | $\downarrow$<br>? | $\downarrow$<br>? | $\downarrow$<br>0<br>:<br>0         | $\downarrow$<br>0 |
| acc.   | ?                 | ?                 | 0<br>:<br>0                         | 0                 |
| $M =$  |                   |                   |                                     |                   |
| t      | 0 ... 0           | ? ... ?           | 0                                   | 1 ... 1           |
| owned  | 0                 | 0                 | $\frac{1}{K}$<br>:<br>$\frac{1}{K}$ | 0                 |
|        | inaccessible      | accessible        | t                                   | owned             |



ranking due to accessible pages

$$r = c M r + (1 - c) q$$

$$r_t = a + cK r_o + (1 - c)q_o$$

$$r_o = c \frac{1}{K} r_t + (1 - c)q_o$$

# Solution

## analysis recap

ranking due to **accessible** pages

teleportation value to pages owned by the spammer

$$r_t = \frac{a}{1 - c^2} + \frac{cK + 1}{1 + c} q_o$$

scaling factor ( $\approx 3.6$ )

**spam** factor (can be made as large as desired)

## solution

teleport only to **trusted** pages (i.e., set  $q_o = 0$ )

can also be used as a method to **identify** spam farms

# PageRank in signed networks

Jung, Jim, Sael, Kang, “Personalized ranking in signed networks using signed random walk with restart,” 2016

<https://ieeexplore.ieee.org/iel7/7837023/7837813/07837935.pdf>

use institutional Sign In with your unipd credentials



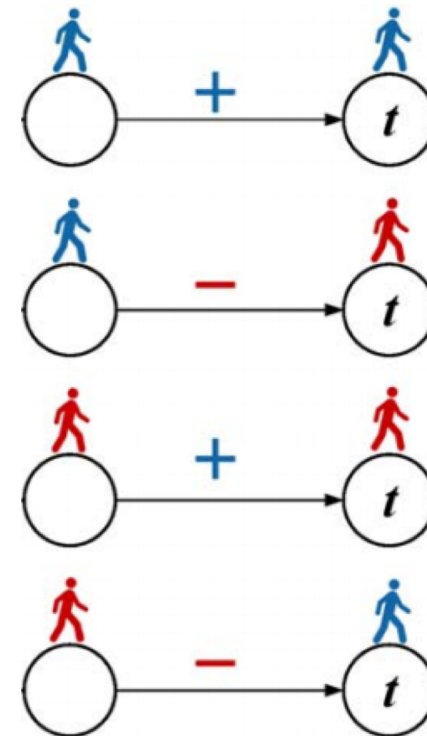
# PageRank in signed networks

Idea:

- Identify + (favourable) and - (adversarial) paths, i.e., ranking values  $r_+$  and  $r_-$  for positive and negative surfers
- Extract positive  $A_+$  and negative  $A_-$  contributions to  $A = A_+ - A_-$
- Normalize the absolute value, to get  $M_+$  and  $M_-$  (with normalized  $M_+ + M_- = 1$ )
- Run a signed random walk

$$r_+ = c M_+ r_+ + c M_- r_- + (1-c) q$$

$$r_- = c M_- r_+ + c M_+ r_-$$



# PageRank in signed networks

signed random walk equation

signed centrality outcome

$$\boxed{r_+ - r_-}$$

damping factor

(column) normalized adjacency matrix

$$r = c M r + (1-c) q_0$$

PageRank vector (centrality)

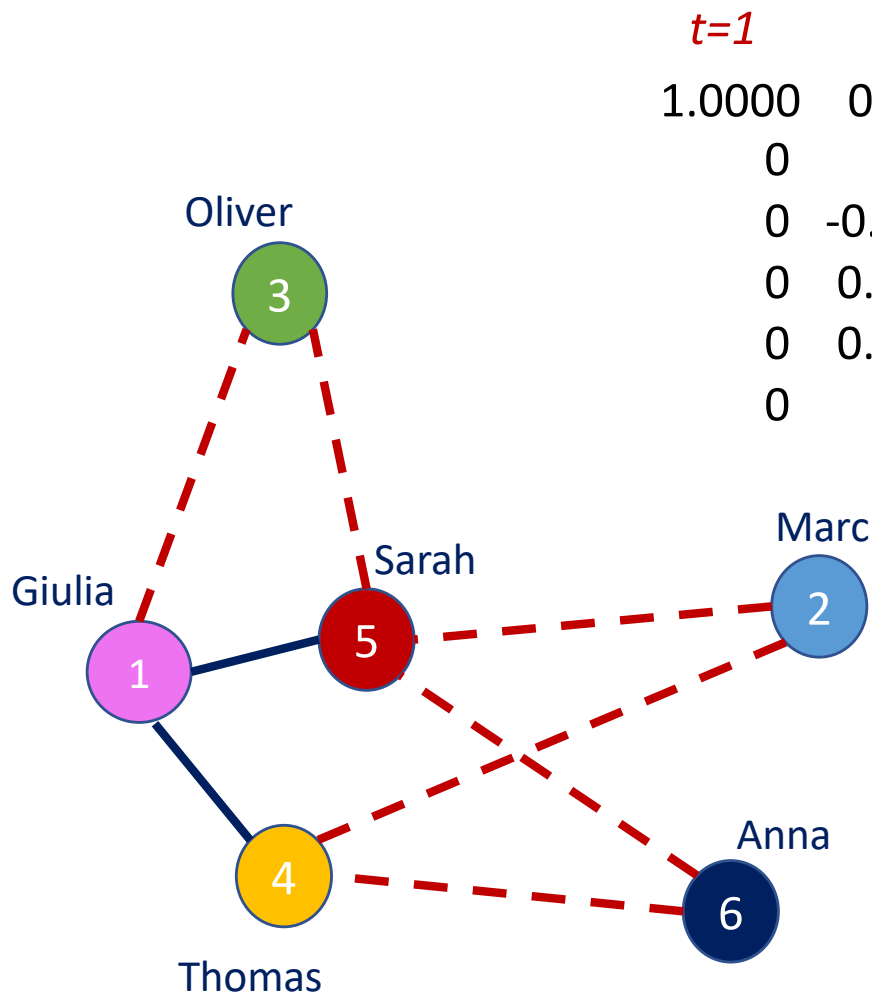
$$r = \begin{bmatrix} r_+ \\ r_- \end{bmatrix}$$

teleportation vector

$$q_0 = \begin{bmatrix} q \\ 0 \end{bmatrix}$$

$$M = \begin{bmatrix} M_+ & M_- \\ M_- & M_+ \end{bmatrix}$$

# Example



|        | <i>t=1</i> | <i>2</i> | <i>3</i> | <i>4</i> | <i>5</i> |
|--------|------------|----------|----------|----------|----------|
| Giulia | 1.0000     | 0.1500   | 0.4109   | 0.2403   | 0.3404   |
| Oliver | 0          | 0        | -0.1405  | -0.0467  | -0.1262  |
| Sarah  | 0          | -0.2833  | -0.1027  | -0.1510  | -0.1275  |
| Marc   | 0          | 0.2833   | 0.0425   | 0.2358   | 0.1078   |
| Thomas | 0          | 0.2833   | 0.1629   | 0.2795   | 0.1719   |
| Anna   | 0          | 0        | -0.1405  | -0.0467  | -0.1262  |

|        | <i>10</i> | <i>20</i> | <i>50</i>     | <i>75</i>     | <i>100</i>    |        |
|--------|-----------|-----------|---------------|---------------|---------------|--------|
| Giulia | 0.2909    | 0.2985    | <b>0.2989</b> | <b>0.2989</b> | <b>0.2989</b> | Giulia |
| Marc   | -0.0848   | -0.0926   | -0.0931       | -0.0931       | -0.0931       | Marc   |
| Oliver | -0.1309   | -0.1313   | -0.1314       | -0.1314       | -0.1314       | Oliver |
| Thomas | 0.1763    | 0.1645    | 0.1638        | 0.1638        | 0.1638        | Thomas |
| Sarah  | 0.2324    | 0.2204    | <b>0.2197</b> | <b>0.2197</b> | <b>0.2197</b> | Sarah  |
| Anna   | -0.0848   | -0.0926   | -0.0931       | -0.0931       | -0.0931       | Anna   |

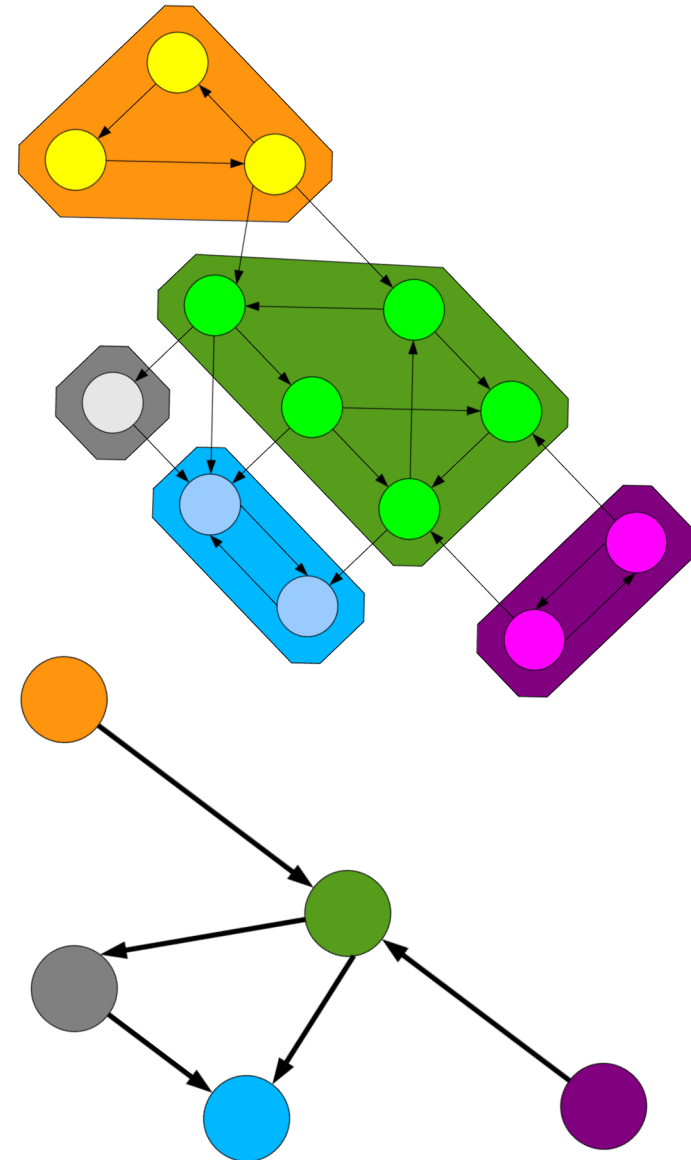
# PageRank eigenstructure

Haveliwala and Kamvar, “The second eigenvalue of the Google matrix,” 2003

<http://ilpubs.stanford.edu:8090/582/1/2003-20.pdf>

# Condensation graph

- Strong connectivity induces a **partition** in disjoint **strongly connected** sets  $\mathcal{V}_1, \mathcal{V}_2, \dots, \mathcal{V}_K$
- By reinterpreting the sets as nodes we obtain a **condensation graph**  $\mathcal{G}^*$  where  $i \rightarrow j$  is an edge if a connection exists between sets  $\mathcal{V}_i \rightarrow \mathcal{V}_j$



# Properties of $\mathcal{G}^*$

- $\mathcal{G}^*$  does not contain **cycles**

otherwise the sets in the cycle would be strongly connected

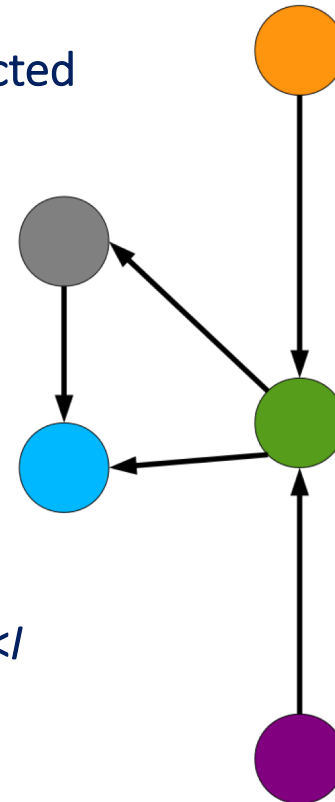
- $\mathcal{G}^*$  has at least one **root** and one **leaf**

and every node in the graph can be reached from one of the roots

- $\mathcal{G}^*$  allows a particular **reordering**

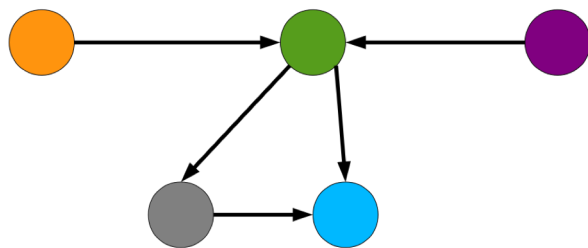
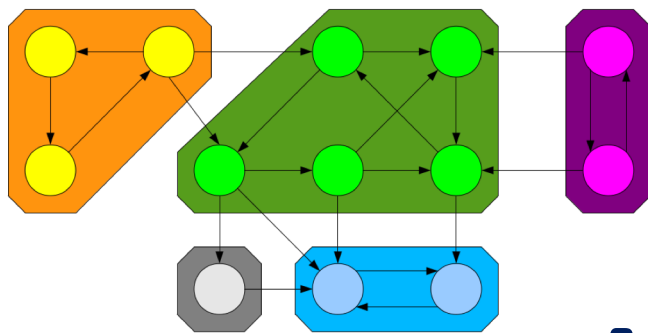
where node  $n_i$  does not reach any of the nodes  $n_j$  with  $j < i$

procedure: identify a root  $n_1$  and remove it from the network, then identify a new root; cycle until all nodes have been selected



# Condensation graph

- The **condensation graph** ordering induces a block-lower-triangular matrix structure on the adjacency matrix



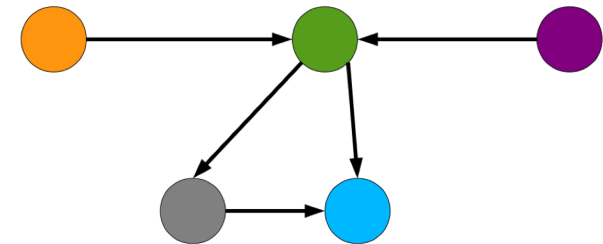
$M =$

|   | 3             | 2             | 5             | 1             | 2 |
|---|---------------|---------------|---------------|---------------|---|
|   | $\frac{1}{3}$ |               |               |               |   |
|   | 1             |               |               |               |   |
| 1 |               |               |               |               |   |
|   |               | $\frac{1}{2}$ |               |               |   |
|   |               | $\frac{1}{2}$ |               |               |   |
|   | $\frac{1}{3}$ |               |               | $\frac{1}{2}$ |   |
|   | $\frac{1}{3}$ | $\frac{1}{2}$ | $\frac{1}{2}$ | $\frac{1}{3}$ |   |
|   |               |               | $\frac{1}{3}$ |               |   |
|   |               | $\frac{1}{2}$ | 1             | $\frac{1}{3}$ |   |
|   |               |               | $\frac{1}{3}$ |               | 0 |
|   |               |               | $\frac{1}{3}$ | $\frac{1}{3}$ | 1 |
|   |               |               |               | $\frac{1}{2}$ | 1 |

blocks in the diagonal are **irreducible** = no block-diagonal form !

# Perron-Frobenius theorem

the eigenvalues of the diagonal blocks, except for the leaves, lie inside the unit circle, i.e.,  $|\lambda| < 1$



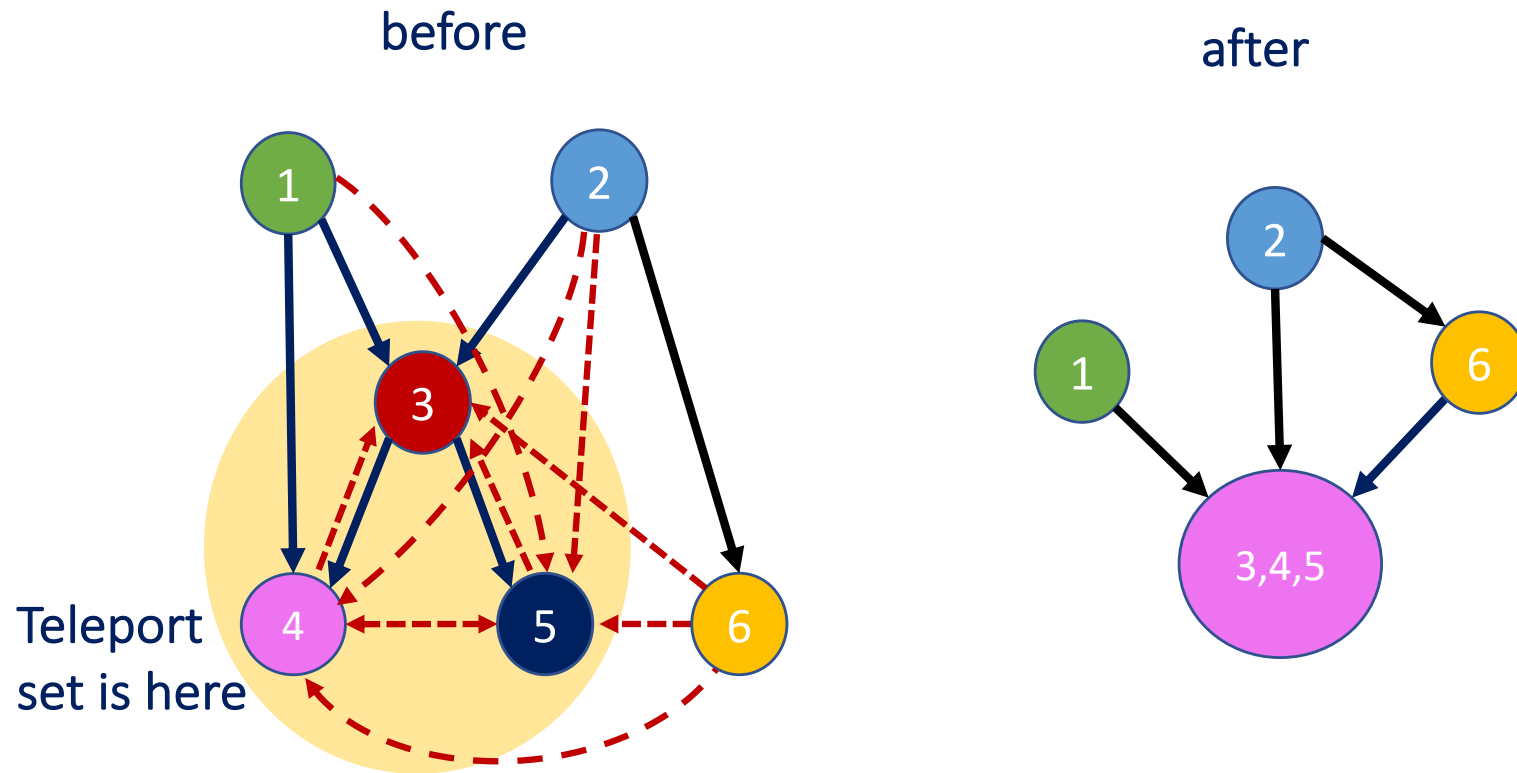
$M =$

|  |               |               |               |               |   |
|--|---------------|---------------|---------------|---------------|---|
|  | 3             | 2             | 5             | 1             | 2 |
|  | $\frac{1}{3}$ |               |               |               |   |
|  | 1             |               |               |               |   |
|  |               | $\frac{1}{2}$ |               |               |   |
|  | $\frac{1}{3}$ | $\frac{1}{2}$ |               | $\frac{1}{2}$ |   |
|  | $\frac{1}{3}$ |               | $\frac{1}{2}$ | $\frac{1}{3}$ |   |
|  |               | $\frac{1}{2}$ | 1             | $\frac{1}{3}$ |   |
|  |               |               | $\frac{1}{3}$ |               | 0 |
|  |               |               | $\frac{1}{3}$ | $\frac{1}{3}$ | 1 |
|  |               |               |               | $\frac{1}{2}$ | 1 |

each leaf-block has at least one eigenvalue in the unit circle,  $\lambda = 1$



# Teleportation ensures one-leaf only



- Hence  $M_1$  carries only **one** eigenvector associated with the eigenvalue  $\lambda=1$

# Lemma

□ PageRank matrix  $M_1 = c M + (1-c) \mathbf{q} \mathbf{1}^T$

□ Normalization property  $\mathbf{1}^T M_1 = \mathbf{1}^T$

□ Jordan form  $M_1 = V J V^{-1}$

carries the right (generalized)  
eigenvectors  $\mathbf{e}_i$  of  $M_1$

carries the  
eigenvalues of  $M_1$

$J =$

$$\begin{bmatrix} \boxed{\begin{matrix} \lambda_1 & 1 \\ & \lambda_1 & 1 \\ & & \lambda_1 \end{matrix}} & & & \\ & \boxed{\begin{matrix} \lambda_2 & 1 \\ & \lambda_2 \end{matrix}} & & \\ & & \boxed{\lambda_3} & \\ & & & \dots \\ & & & & \boxed{\begin{matrix} \lambda_n & 1 \\ & \lambda_n \end{matrix}} \end{bmatrix}$$

$$\begin{aligned} \mathbf{1}^T M_1 V &= \mathbf{1}^T V \\ &= \mathbf{1}^T V J \end{aligned} \quad \Rightarrow \quad \underbrace{\mathbf{1}^T V}_{\rho} \underbrace{(J - I)}_{\text{only one eigenvalue is 0}} = 0$$

□ Hence  $\mathbf{1}^T \mathbf{e}_i = 0$  for  $i > 1$ , i.e., except for the eigenvector associated with eigenvalue 1

# Main result

$$\underbrace{M_1 e_i = c M e_i}_{\text{same eigenvalues of } M, \text{ but multiplied by } c !!!} + (1-c) \cancel{q \mathbf{1}^T e_i} \quad \text{for } i > 1$$

same eigenvalues of  $M$ ,  
but multiplied by  $c$  !!!



- ❑  $M_1$  has **one** eigenvalue equal to **1**
- ❑ The remaining eigenvalues satisfy  $|\lambda| \leq c$

# Approximate PageRank

Andersen, Chung, Lang, “Local graph partitioning using PageRank vectors,” 2006

<https://ieeexplore.ieee.org/stamp/stamp.jsp?arnumber=4031383>

use [institutional Sign In](#) with your unipd credentials

# Approximate algorithm

□ Start from  $\mathbf{u} = \mathbf{0}$  and  $\mathbf{v} = \mathbf{q}$

□ To all the nodes  $i$  satisfying  $\underline{v_i} > \varepsilon \underline{d_i}/D$  apply the push operation

precision  
degree of node  $i$   
sum of the degrees

$u$  constantly increases  $\longrightarrow \mathbf{u}^+ = \mathbf{u} + (1-c) \delta$   $\longleftarrow$  only one active element in position  $i$  with value  $v_i$

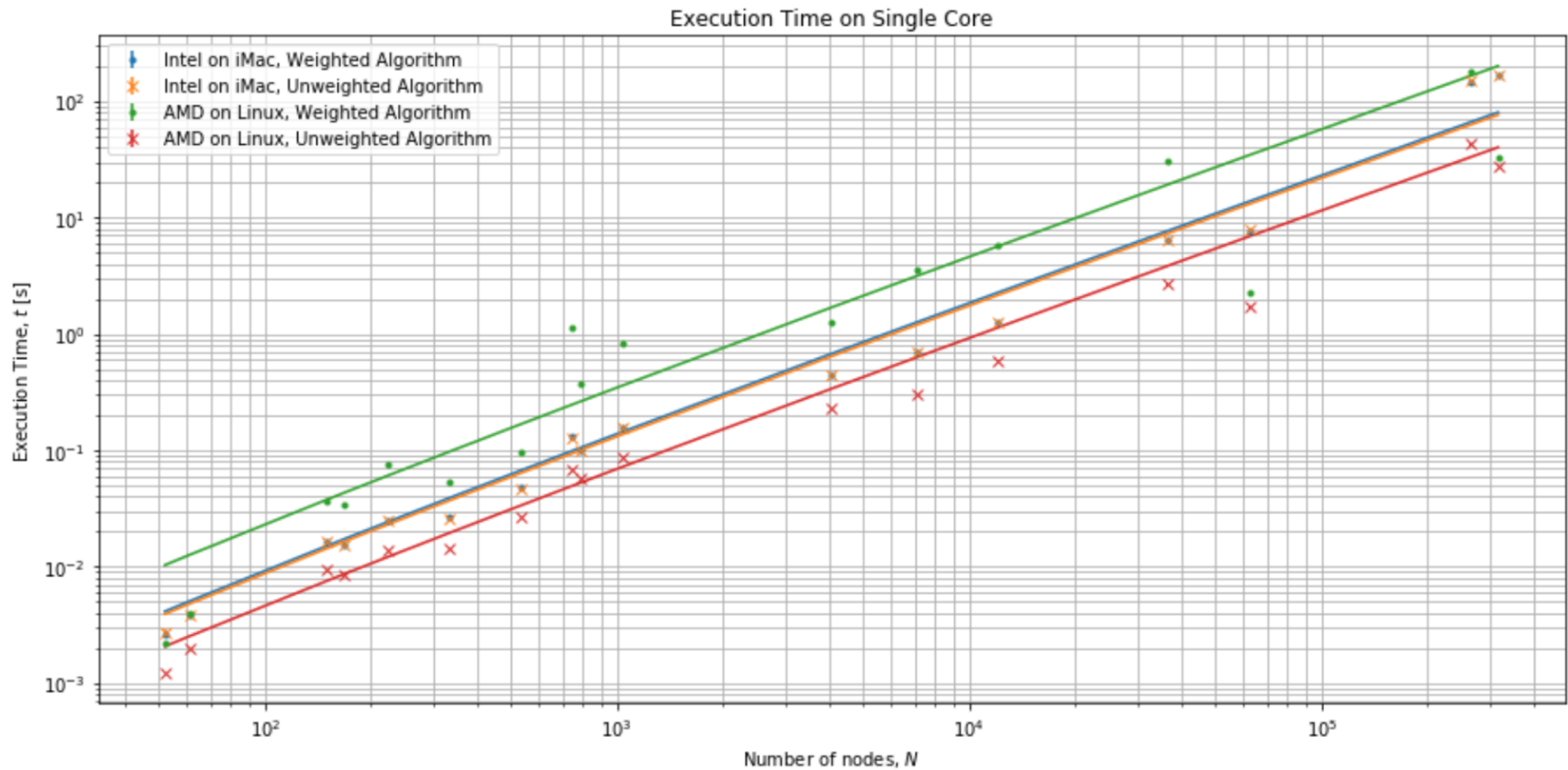
$v$  always positive  $\longrightarrow \mathbf{v}^+ = \mathbf{v} - \delta + c \mathbf{M} \delta$

□ Returns  $\mathbf{r} \simeq \mathbf{u}$  with precision  $|\mathbf{r} - \mathbf{u}|_1 < \varepsilon$

□ It is simple

# Scalability properties

(Francesco Barbato & Tommaso Boccato, 2020)



# PageRank linearity

column stochastic matrix  $\mathbf{1}^T \mathbf{M} = \mathbf{1}^T$

□ PageRank equation  $\mathbf{r}_q = c \mathbf{M} \mathbf{r}_q + (1-c) \mathbf{q}$

stochastic ranking vector  
 $\mathbf{1}^T \mathbf{r}_q = 1, \mathbf{r}_q \geq 0$

stochastic Teleport vector  
 $\mathbf{1}^T \mathbf{q} = 1$

□ Alternative equation  $\mathbf{r}_q = (\mathbf{I} - c \mathbf{M})^{-1} (1-c) \mathbf{q}$

linear in  $\mathbf{q}$

$$\mathbf{r}_{a\mathbf{u}+b\mathbf{v}} = a \mathbf{r}_u + b \mathbf{r}_v$$

# Alternative equation

one-step random walk

□ PageRank equation  $r_q = c r_{Mq} + (1-c) q$



□  $r_q = (I - c M)^{-1} (1-c) q$

□  $r_q = (1-c) \sum (c M)^k q$

□  $M r_q = (1-c) \sum (c M)^k M q$

□  $M r_q = r_{Mq}$



# Main property of push: $r_q = u + r_v$

- At starting point  $u = 0$  and  $v = q$  imply  $r_q = 0 + r_q$
- The following steps are proved by induction

$$u^+ = u + (1-c) \delta$$

$$v^+ = v - \delta + c M \delta$$



$$u^+ + r_{v^+} = u + (1-c) \delta + \overbrace{r_v - r_\delta + c r_{M\delta}}^{\text{by linearity}}$$



$$u^+ + r_{v^+} = u + r_v = r_q$$

# Precision guarantee $\|r_q - u\|_1 < \varepsilon$

- The push property implies  $r_q = u + r_v$
- Hence  $\|r_q - u\|_1 = \|r_v\|_1 = \mathbf{1}^T r_v$
- The PageRank equation is  $r_v = c M r_v + (1-c) v$
- Hence  $\mathbf{1}^T r_v = c \underbrace{\mathbf{1}^T M}_{\mathbf{1}^T} r_v + (1-c) \mathbf{1}^T v$  so that  $\mathbf{1}^T r_v = \mathbf{1}^T v$
- As a result  $\|r_q - u\|_1 = \mathbf{1}^T v < \sum \varepsilon d_i / D = \varepsilon$

# Lazy PageRank

□ Lazy PageRank  $r = a \underline{M_2} r + (1-a) q$

$$M_2 = b I + (1-b) M$$

□ **Lazy** because a fraction  $b$  of the times the surfer stays where she/he is

□ Equivalent to  $r = c \underline{M} r + (1-c) q$

$$c = a(1-b) / \underbrace{(1-ab)}_{\text{slower algorithm}} < a$$

slower algorithm

# Lessons learned

- ❑ Importance of Teleport vector
- ❑ PageRank can measure similarity
- ❑ PageRank can be extended to signed networks (with a trick)
- ❑ Reliable and scalable implementations exist

