Introduction
○○○
Related Work
○
JXP Algorithm
○○○○○○○○○○
Mathematical Analysis
○○○○
Experimental Results
○○○○○
Conclusion
○

# Efficient and Decentralized PageRank Approximation in a P2P Network

Josiane Xavier Parreira, Debora Donato, Sebastian Michel, Gerhard Weikum

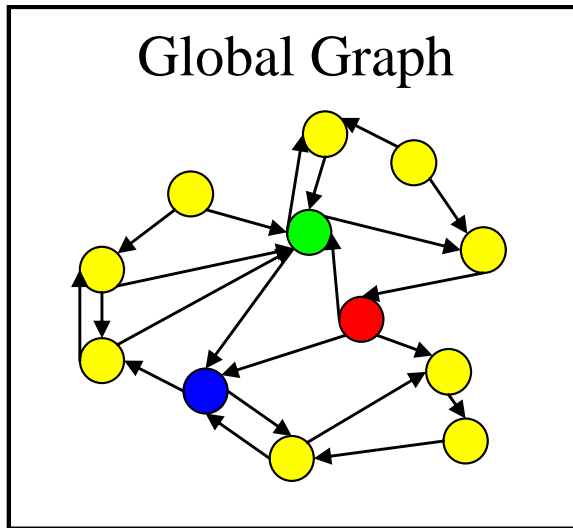Max Planck Institute for Computer Science, Universita di Roma

May 20$^{th}$, 2006

Introduction
○○○
Related Work
○
JXP Algorithm
○○○○○○○○○○
Mathematical Analysis
○○○○
Experimental Results
○○○○○
Conclusion
○

# Outline

Introduction | Related Work | JXP Algorithm | Mathematical Analysis | Experimental Results | Conclusion
●○○ | ○ | ○○○○○○○○○○ | ○○○○ | ○○○○○ | ○

Introduction

# Motivation

Scenario: A graph...

Introduction | Related Work | JXP Algorithm | Mathematical Analysis | Experimental Results | Conclusion
○●○ | ○ | ○○○○○○○○○○ | ○○○○ | ○○○○○ | ○

Introduction
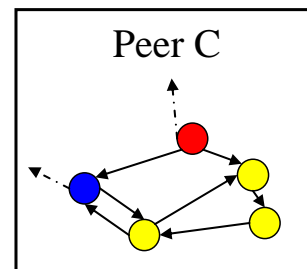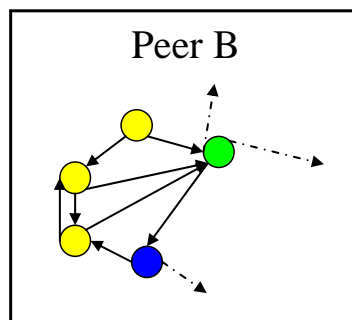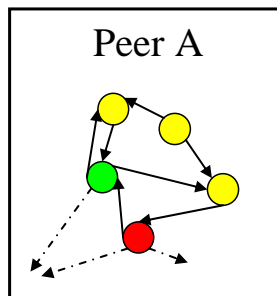
# Motivation

Scenario: A graph...



...distributed into peers in a P2P network.

# Motivation

## Goal

Compute "global" authority scores of pages in the network.

---

# Motivation

## Goal

Compute "global" authority scores of pages in the network.

## Problems

- Peers have only local (incomplete) information
- Pages might link to or be linked by pages at other peers
- No control over overlaps between local graphs

# Related Work

## Efficient PR

- Graph Aggregation [Broder et al., WWW'04]
- Iterative Aggregation [Langville & Meyer, WWW'04]

## Decentralized PR

- *Local PageRank & ServerRank* [Wang & DeWitt, VLDB'04]
- *BlockRank* [Kamvar et al., Stanford Tech. Report'03]

## Markov Chains Aggregation/Disaggregation Techniques

- Kemeny & Snell [1963]
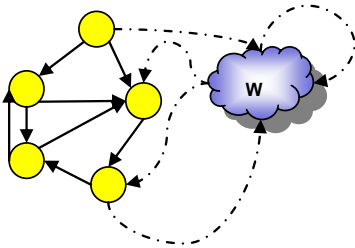- Stewart [1994]
- Meyer [2000]

# JXP Algorithm

## JXP Algorithm

- Decentralized algorithm for computing authority scores of pages in a P2P Network
- Runs locally at every peer
- Combines local PageRank computations + Meetings between peers
- JXP scores converge to the true global PageRank scores

| Introduction | Related Work | JXP Algorithm | Mathematical Analysis | Experimental Results | Conclusion |
|---|---|---|---|---|---|
| ○○○ | ○ | ○●○○○○○○○○ | ○○○○ | ○○○○○ | ○ |

World node

# World Node



---

| Introduction | Related Work | JXP Algorithm | Mathematical Analysis | Experimental Results | Conclusion |
|---|---|---|---|---|---|
| ○○○ | ○ | ○●○○○○○○○○ | ○○○○ | ○○○○○ | ○ |

World node

# World Node

- Special node added to each local graph

| Introduction | Related Work | JXP Algorithm | Mathematical Analysis | Experimental Results | Conclusion |
|---|---|---|---|---|---|
| ○○○ | ○ | ○●○○○○○○○○ | ○○○○ | ○○○○○ | ○ |

World node

# World Node



- Special node added to each local graph
- Represents all pages in the network that do not belong to local graph

| Introduction | Related Work | JXP Algorithm | Mathematical Analysis | Experimental Results | Conclusion |
|---|---|---|---|---|---|
| ○○○ | ○ | ○●○○○○○○○○ | ○○○○ | ○○○○○ | ○ |

World node

# World Node



- Special node added to each local graph
- Represents all pages in the network that do not belong to local graph
- "Special features":

| Introduction | Related Work | JXP Algorithm | Mathematical Analysis | Experimental Results | Conclusion |
|---|---|---|---|---|---|
| ○○○ | ○ | ○●○○○○○○○○ | ○○○○ | ○○○○○ | ○ |

World node

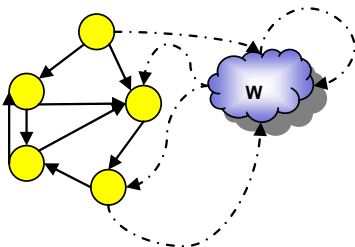# World Node



- Special node added to each local graph
- Represents all pages in the network that do not belong to local graph
- "Special features":
  - All links from local pages to external pages point to World Node

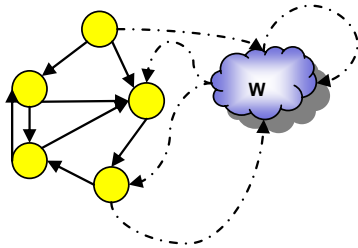| Introduction | Related Work | JXP Algorithm | Mathematical Analysis | Experimental Results | Conclusion |
|---|---|---|---|---|---|
| ○○○ | ○ | ○●○○○○○○○○ | ○○○○ | ○○○○○ | ○ |

World node

# World Node



- Special node added to each local graph
- Represents all pages in the network that do not belong to local graph
- "Special features":
  - All links from local pages to external pages point to World Node
  - External links (discovered during meetings) are represented at the World Node

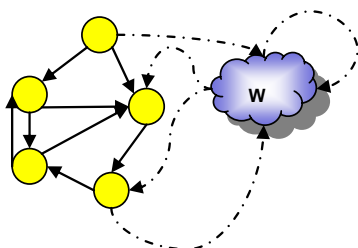| Introduction | Related Work | JXP Algorithm | Mathematical Analysis | Experimental Results | Conclusion |
| --- | --- | --- | --- | --- | --- |
| ○○○ | ○ | ○●○○○○○○○○ | ○○○○ | ○○○○○ | ○ |

World node

# World Node



- Special node added to each local graph
- Represents all pages in the network that do not belong to local graph
- "Special features":
  - All links from local pages to external pages point to World Node
  - External links (discovered during meetings) are represented at the World Node
  - Transitions prob. from external incoming links are stored; World Node outgoing links are weighted to reflect score mass given by original link

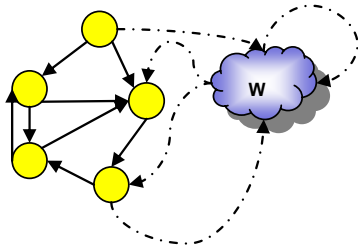| Introduction | Related Work | JXP Algorithm | Mathematical Analysis | Experimental Results | Conclusion |
| --- | --- | --- | --- | --- | --- |
| ○○○ | ○ | ○●○○○○○○○○ | ○○○○ | ○○○○○ | ○ |

World node

# World Node



- Special node added to each local graph
- Represents all pages in the network that do not belong to local graph
- "Special features":
  - All links from local pages to external pages point to World Node
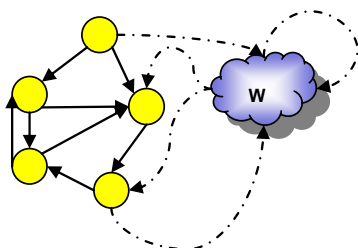  - External links (discovered during meetings) are represented at the World Node
  - Transitions prob. from external incoming links are stored; World Node outgoing links are weighted to reflect score mass given by original link
  - Self-loop link to represent transitions among external pages

Introduction · · · · Related Work · JXP Algorithm · · · · · · · · · · Mathematical Analysis · · · · Experimental Results · · · · · Conclusion ·

JXP Algorithm

# The Algorithm

## Initialization step

- Local graph is extended by adding the world node;
- PageRank is computed in the extended graph → JXP scores

Introduction · · · · Related Work · JXP Algorithm · · · · · · · · · · Mathematical Analysis · · · · Experimental Results · · · · · Conclusion ·

JXP Algorithm

# The Algorithm

## Initialization step

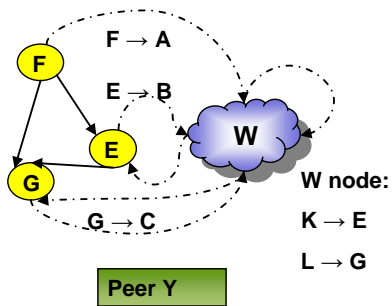- Local graph is extended by adding the world node;
- PageRank is computed in the extended graph → JXP scores

## Main Algorithm (for every $P_i$ in the network)

- Select $P_j$ to meet
- Update world node
  - Add links from $P_j$ that point to pages in $P_i$
  - If a link already exists at the world node, the score of the source page is updated by taking the highest of both scores
- Compute PageRank → JXP scores

| Introduction | Related Work | JXP Algorithm | Mathematical Analysis | Experimental Results | Conclusion |
|---|---|---|---|---|---|
| ○○○ | ○ | ○○○●○○○○○ | ○○○○ | ○○○○○ | ○ |

JXP Algorithm

# Example

| Introduction | Related Work | JXP Algorithm | Mathematical Analysis | Experimental Results | Conclusion |
|---|---|---|---|---|---|
| ○○○ | ○ | ○○○●○○○○○ | ○○○○ | ○○○○○ | ○ |

JXP Algorithm

# Example

| Introduction | Related Work | JXP Algorithm | Mathematical Analysis | Experimental Results | Conclusion |
|---|---|---|---|---|---|
| ○○○ | ○ | ○○○●○○○○○ | ○○○○ | ○○○○○ | ○ |

JXP Algorithm

# Example

| Introduction | Related Work | JXP Algorithm | Mathematical Analysis | Experimental Results | Conclusion |
|---|---|---|---|---|---|
| ○○○ | ○ | ○○○○●○○○○○ | ○○○○ | ○○○○○ | ○ |

Peer Selection Strategy

# Peer Selection Strategy

## Motivation

- Peers' contribution for the convergence are different
- Finding most promising peers would speed up convergence
- "Quality indicator": Amount of outgoing links of a peer in the network that are also incoming links in the local graph

| Introduction | Related Work | JXP Algorithm | Mathematical Analysis | Experimental Results | Conclusion |
| ○○○ | ○ | ○○○○●○○○○○ | ○○○○ | ○○○○○ | ○ |

Peer Selection Strategy

# Peer Selection Strategy

## Motivation

- Peers' contribution for the convergence are different
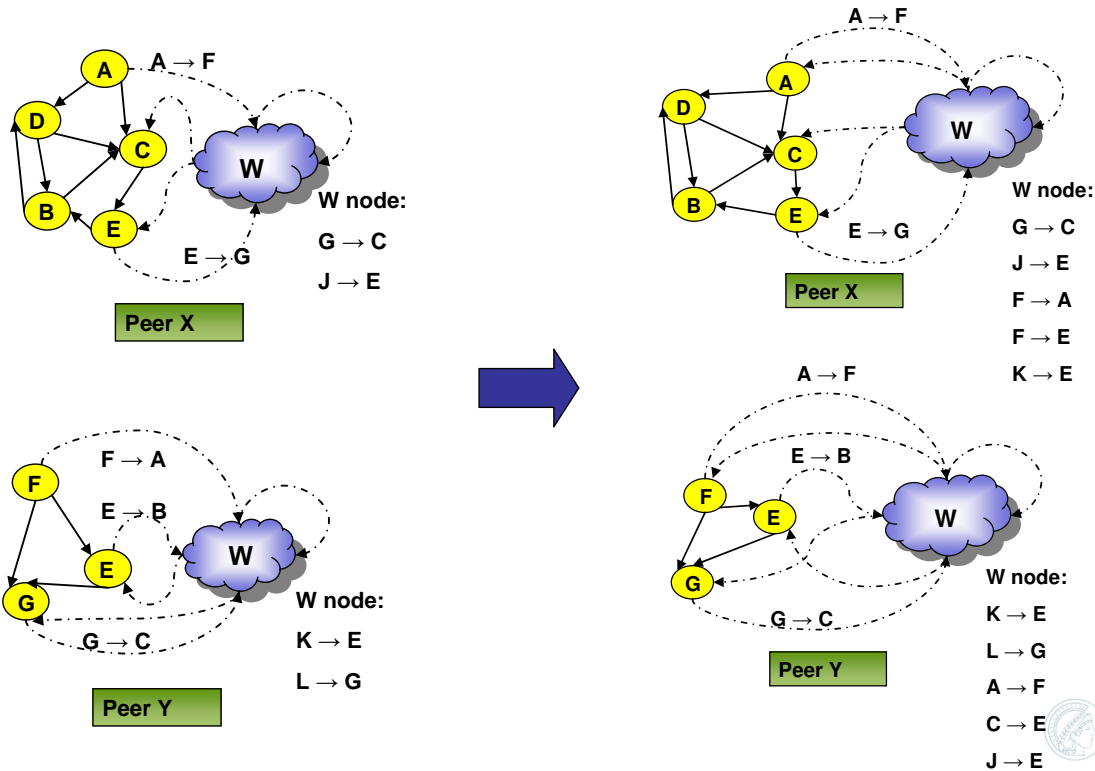- Finding most promising peers would speed up convergence
- "Quality indicator": Amount of outgoing links of a peer in the network that are also incoming links in the local graph

| Introduction | Related Work | JXP Algorithm | Mathematical Analysis | Experimental Results | Conclusion |
| ○○○ | ○ | ○○○○○●○○○○ | ○○○○ | ○○○○○ | ○ |

Peer Selection Strategy

# Peer Selection Strategy

## Good strategy

Find promising peers without increasing much bandwidth consumption

- Caching + statistical synopses

# Peer Selection Strategy

## Good strategy

Find promising peers without increasing much bandwidth consumption

- Caching + statistical synopses

## Statistical synopses

Approximation technique for comparing data of different peers without explicitly transferring their contents.

- Compact representation of sets
- Can be used to estimate cardinality of the intersection between two sets
- JXP uses Min-Wise Independent Permutations (MIPs)

# MIPs

## MIPs [Broder et al., 2000]

- Computes $N$ permutations of the set's elements
- Permutation uses linear hash function
  $h_i(x) := a_i * x + b_i \ mod \ U$
- For each of $N$ permutations, it determines the minimum hash value and stores it in a $N$-dimensional vector
- Idea: Each element has same probability of becoming the minimum element
- MIPs vector are then used to estimate the intersection between two sets

| Introduction | Related Work | JXP Algorithm | Mathematical Analysis | Experimental Results | Conclusion |
| ○○○ | ○ | ○○○○○○○●○○ | ○○○○ | ○○○○○ | ○ |

Peer Selection Strategy

## MIPs Example



*pageID* set

| 17 | 21 | 3 | 12 | 24 | 8 |

$h_1 = (7x + 3) \bmod 51$

| 20 | 48 | 24 | 36 | 18 | 8 |

$h_2 = (5x + 6) \bmod 51$

| 40 | 9 | 21 | 15 | 24 | 46 |

$h_N = (3x + 9) \bmod 51$

| 9 | 21 | 18 | 45 | 30 | 33 |

Apply Permutations to all *pageID*'s

| 8 |
| 9 |
| ⋮ |
| 9 |

Create MIP Vector from Minima of Permutations

N

Figure: Example of Min-Wise Permutations.

| Introduction | Related Work | JXP Algorithm | Mathematical Analysis | Experimental Results | Conclusion |
| ○○○ | ○ | ○○○○○○○○●○ | ○○○○ | ○○○○○ | ○ |

Peer Selection Strategy

## Pre-meetings Strategy

- Each peer $P_i$ computes $local(P_i)$ and $successors(P_i)$ MIPs vectors
- When $P_i$ meets $P_j$
  - Uses MIPs vectors to estimate percentage of local pages pointed by pages in $P_j$
  - If percentage above threshold, $P_i$ caches $P_j$'s ID
  - Uses MIPs again to estimate overlap between overlap between the two local graphs
  - If there is high overlap, peers exchange their list of cached ID's and store them in a temporary list
  - Idea: Peers on the temporary list are potential candidates for the next meeting

Introduction    Related Work    JXP Algorithm    Mathematical Analysis    Experimental Results    Conclusion
○○○             ○               ○○○○○○○○○●         ○○○○                   ○○○○○                  ○

Peer Selection Strategy

# Pre-meetings Strategy

### Pre-meetings phase

- $P_j$ contacts peers on the temporary list and ask for their MIPs vectors
- Assign scores to each peer
- For next (real) meeting, $P_i$ chooses $P_k$ where
  - $P_k$ is best scored peer in temporary list, with prob. $\alpha$
  - $P_k$ is one of the already cached peers, with prob. $\beta$
  - $P_k$ is a random peer, with prob. $(1 - \alpha - \beta)$

# Mathematical Analysis

### Assumptions

Global transition matrix $\mathbf{C}_{N \times N}$ and global stationary distribution vector $\boldsymbol{\pi}$

### Local transition matrix and local stationary distr. (JXP scores)

$$\mathbf{P} = \begin{pmatrix} p_{11} & \cdots & p_{1n} & p_{1w} \\ \vdots & \cdots & \vdots & \vdots \\ p_{n1} & \cdots & p_{nn} & p_{nw} \\ \hline p_{w1} & \cdots & p_{wn} & p_{ww} \end{pmatrix}$$

$$p_{ij} = \begin{cases} \frac{1}{out(i)} & \text{if } \exists\; i \to j \\ 0 & \text{otherwise} \end{cases}$$

$$p_{iw} = \sum_{\substack{i \to r \\ r \notin G}} \frac{1}{out(i)}$$

$$\boldsymbol{\alpha} = \begin{pmatrix} \alpha_1 & \cdots & \alpha_n & | & \alpha_w \end{pmatrix}^T$$

for every $i, j$, $1 \le i, j \le n$. ($G$ is the set of local pages)

Introduction
ooo

Related Work
o

JXP Algorithm
oooooooooo

Mathematical Analysis
o●oo

Experimental Results
ooooo

Conclusion
o

# Mathematical Analysis

## World Node transitions prob.

$$p_{wi}^t = \sum_{\substack{r \to i \\ r \in W^t}} \frac{\alpha(r)^t}{out(i)} \cdot \frac{1}{\alpha_w^{t-1}} \qquad p_{ww}^t = 1 - \sum_{i=1}^n p_{wi}^t$$

$W^t$: Set of pages represented at the World Node during meeting $t$

---

Introduction
ooo

Related Work
o

JXP Algorithm
oooooooooo

Mathematical Analysis
o●oo

Experimental Results
ooooo

Conclusion
o

# Mathematical Analysis

## World Node transitions prob.

$$p_{wi}^t = \sum_{\substack{r \to i \\ r \in W^t}} \frac{\alpha(r)^t}{out(i)} \cdot \frac{1}{\alpha_w^{t-1}} \qquad p_{ww}^t = 1 - \sum_{i=1}^n p_{wi}^t$$

$W^t$: Set of pages represented at the World Node during meeting $t$

## Random Jumps

$$\mathbf{P}' = \epsilon\, \mathbf{P} + (1 - \epsilon)\frac{1}{N} \left(\; 1 \quad \dots \quad 1 \;\middle|\; (N - n)\; \right)$$

# Mathematical Analysis

## Meeting Step

Considering one link addition/update at a time

$$\mathbf{P}^t = \mathbf{P}^{t-1} + \mathbf{E} \quad \mathbf{E} = \begin{pmatrix} 0 & & & & \dots & & & 0 & 0 \\ \vdots & & & & \dots & & & \vdots & \vdots \\ 0 & & & & \dots & & & 0 & 0 \\ \hline 0 & \dots & 0 & \delta & 0 & \dots & 0 & -\delta \end{pmatrix}$$

# Mathematical Analysis

## Meeting Step

Considering one link addition/update at a time

$$\mathbf{P}^t = \mathbf{P}^{t-1} + \mathbf{E} \quad \mathbf{E} = \begin{pmatrix} 0 & & & & \dots & & & 0 & 0 \\ \vdots & & & & \dots & & & \vdots & \vdots \\ 0 & & & & \dots & & & 0 & 0 \\ \hline 0 & \dots & 0 & \delta & 0 & \dots & 0 & -\delta \end{pmatrix}$$

## Theorem 1

*The JXP score of the world node, at every peer in the network, is monotonically non-increasing.*

# Mathematical Analysis

## Theorem 2

*Consider the true stationary probabilities (PR scores) of pages $i \in G$ and the world node $w$, $\pi_i$ and $\pi_w$, and their JXP scores after $t$ meetings $\alpha_i^t$ and $\alpha_w^t$. The following holds throughout all JXP meetings:*
*$0 < \alpha_i^t \leq \pi_i$ for $i \in G$ and $\pi_w \leq \alpha_w^t < 1$.*

# Mathematical Analysis

## Theorem 2

*Consider the true stationary probabilities (PR scores) of pages $i \in G$ and the world node $w$, $\pi_i$ and $\pi_w$, and their JXP scores after $t$ meetings $\alpha_i^t$ and $\alpha_w^t$. The following holds throughout all JXP meetings:*
*$0 < \alpha_i^t \leq \pi_i$ for $i \in G$ and $\pi_w \leq \alpha_w^t < 1$.*

## Theorem 3

*In a fair series of JXP meetings, the JXP scores of all nodes converge to the true global PR scores.*

| Introduction | Related Work | JXP Algorithm | Mathematical Analysis | Experimental Results | Conclusion |
|---|---|---|---|---|---|
| ○○○ | ○ | ○○○○○○○○○○ | ○○○○ | ●○○○○ | ○ |

JXP Accuracy and Convergence

# Setup

| Introduction | Related Work | JXP Algorithm | Mathematical Analysis | Experimental Results | Conclusion |
|---|---|---|---|---|---|
| ○○○ | ○ | ○○○○○○○○○○ | ○○○○ | ●○○○○ | ○ |

JXP Accuracy and Convergence

# Setup

## Amazon collection
- 55,196 pages
- 237,160 links
- 10 categories

## Web collection
- 103,591 pages
- 1,633,276 links
- 10 categories

| Introduction | Related Work | JXP Algorithm | Mathematical Analysis | Experimental Results | Conclusion |
|---|---|---|---|---|---|
| ○○○ | ○ | ○○○○○○○○○○ | ○○○○ | ●○○○○ | ○ |

JXP Accuracy and Convergence

# Setup

## Amazon collection
- 55,196 pages
- 237,160 links
- 10 categories

## Web collection
- 103,591 pages
- 1,633,276 links
- 10 categories
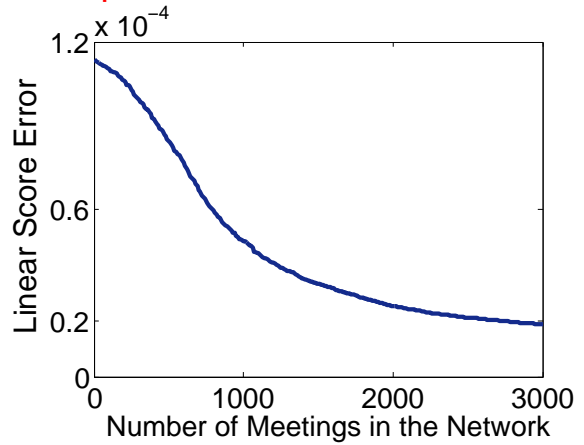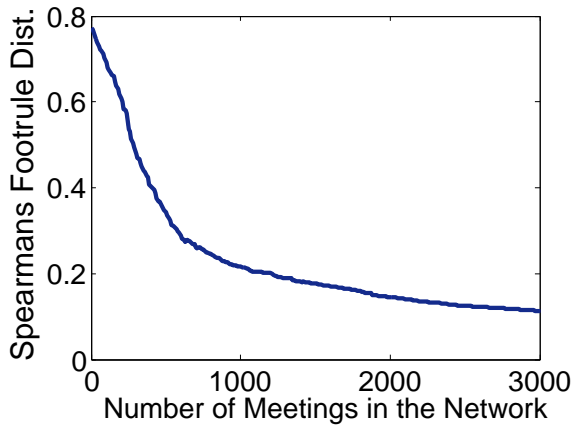
## Setup
- 100 peers
- 10 peers/category

---

## Evaluation Measures
- "Global" JXP ranking vs. Global PageRank ranking
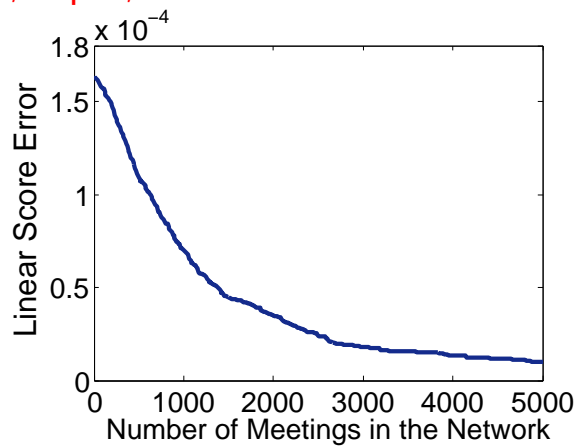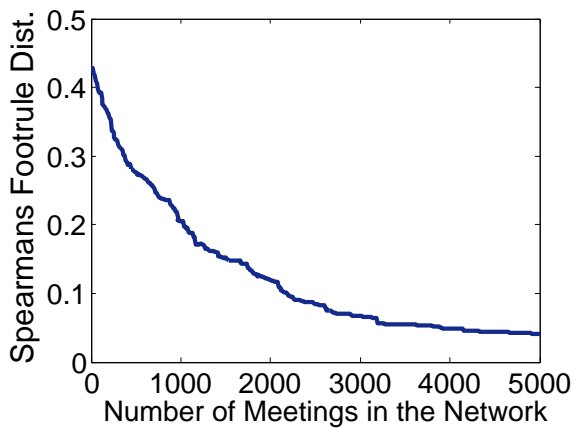- Spearman's Footrule Distance at top-k
- Linear Score Error at top-k

Introduction
○○○
Related Work
○
JXP Algorithm
○○○○○○○○○○
Mathematical Analysis
○○○○
Experimental Results
○●○○○
Conclusion
○

JXP Accuracy and Convergence

## Overall performance comparison



Amazon Collection, top-1,000

Introduction
○○○
Related Work
○
JXP Algorithm
○○○○○○○○○○
Mathematical Analysis
○○○○
Experimental Results
○○●○○
Conclusion
○

JXP Accuracy and Convergence

## Overall performance comparison



Web Collection, top-1,000

# JXP in P2P Search

# JXP in P2P Search

JXP integrated into the Minerva search engine for P2P networks.
(Minerva Project Website: http://www.minerva-project.org)

## Setup

- Bigger set of Web Collection (250,760 docs and 3,123,993 links)

- 40 peers, high overlap

- 15 queries, using the Minerva query routing mechanism

- Results were ranked in two ways:
  - tf*idf only
  - weighted sum of tf*idf and JXP scores

- Precision at top-10 measured

Introduction | Related Work | JXP Algorithm | Mathematical Analysis | **Experimental Results** | Conclusion
○○○ | ○ | ○○○○○○○○○○ | ○○○○ | ○○○○○● | ○

JXP in P2P Search

# Results

| Query | tf*idf | (0.6 tf*idf + 0.4 JXP) |
|:---:|:---|:---:|
| affirmative action | 40% | 40% |
| amusement parks | 60% | 60% |
| armstrong | 20% | **80%** |
| basketball | 20% | **60%** |
| blues | 20% | 20% |
| censorship | **30%** | 20% |
| cheese | 40% | **60%** |
| iraq war | **50%** | 30% |
| jordan | 40% | 40% |
| moon landing | **90%** | 70% |
| movies | 30% | **100%** |
| roswell | 30% | **70%** |
| search engines | 20% | **60%** |
| shakespeare | 60% | **80%** |
| table tennis | 50% | **70%** |
| **Average** | 40% | **57%** |

Introduction | Related Work | JXP Algorithm | Mathematical Analysis | Experimental Results | **Conclusion**
○○○ | ○ | ○○○○○○○○○○ | ○○○○ | ○○○○○ | ●

Conclusion

# Conclusions and Ongoing Work

## Conclusions

- JXP algorithm for dynamically computing authority scores of pages distributed in a P2P network
- Combines local PageRank computation with meetings between peers
- JXP scores are proved to converge to global PageRank scores

## Ongoing Work

- Integrate JXP into the query routing mechanism
- JXP in dynamic networks