



RankMass Crawler

A Crawler with High Personalized
PageRank Coverage Guarantee

Junghoo Cho
Uri Schonfeld

Work done in UCLA



Motivation

- Impossible to download the entire web:
 - Example: many pages from one calendar
- When can we stop?
- How to gain the most benefit from the pages we download



Main Issues

- Crawler Guarantee:
 - *guarantee* on how much of the “important” part of the Web they “cover” when they stop crawling
 - If we *don't see the pages*, how do we know how important they are?
- Crawler Efficiency:
 - Download “important” pages early during a crawl
 - Obtain coverage with a *min number of downloads*



Outline

- Formalize coverage metric
- L-Neighbor: Crawling with RankMass guarantee
- RankMass: Crawling to achieve high RankMass
- Windowed RankMass: How greedy do you want to be?
- Experimental Results



Web Coverage Problem

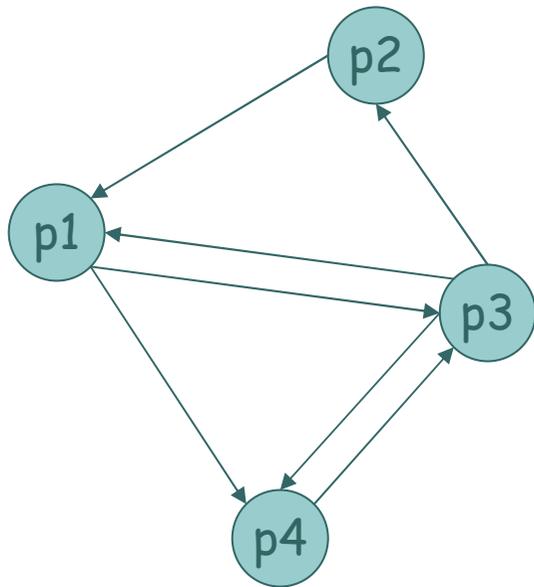
- D - The potentially infinite set of documents of the web
- D_c - The finite set of documents in our document collection
- **Assign** importance **weights** to each page



Web Coverage Problem

- What weights? Per query? Topic? Font?
- PageRank? Why PageRank?
 - Useful as importance measure
 - Random surfer.
 - Effective for ranking.

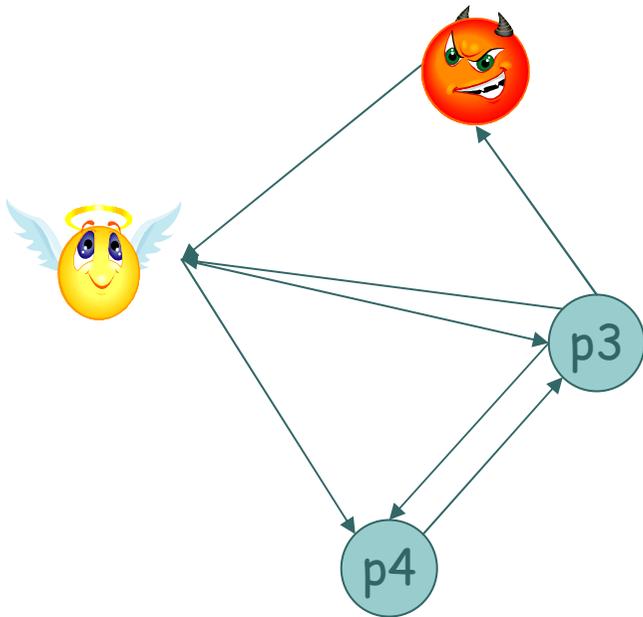
● ● ● | PageRank a Short Review



$$r_i = d \left[\sum_{p_j \in I(p_i)} \frac{r_j}{c_j} \right] + (1-d) \frac{1}{|D|}$$

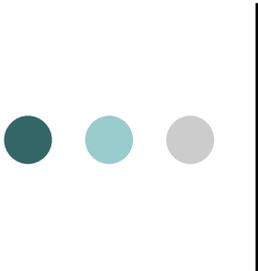
● ● ● | Now it's Personal

- Personal, TrustRank, General



$$T = \begin{bmatrix} t_1 \\ t_2 \\ M \end{bmatrix}$$

$$r_i = d \left[\sum_{p_j \in I(p_i)} \frac{r_j}{c_j} \right] + (1-d)t_i$$



RankMass Defined

- Using personalized pagerank formally define RankMass of D_C :

$$RM(D_C) = \sum_{p_i \in D_C} r_i$$

- Coverage Guarantee:

- We seek a crawler that given ϵ , when it stops the downloaded pages D_C :

$$RM(D_C) = \sum_{p_i \in D_C} r_i > 1 - \epsilon$$

- Efficient crawling:

- We seek a crawler that, for a given N , downloads $|D_C|=N$ s.t. $RM(D_C)$ is greater or equal to any other $|D_C|=N$, $D_C \subseteq D$



How to Calculate RankMass

- Based on PageRank
- How do you compute $RM(Dc)$ without downloading the entire web
- We can't compute the exact but can lower bound
- Let's start a simple case

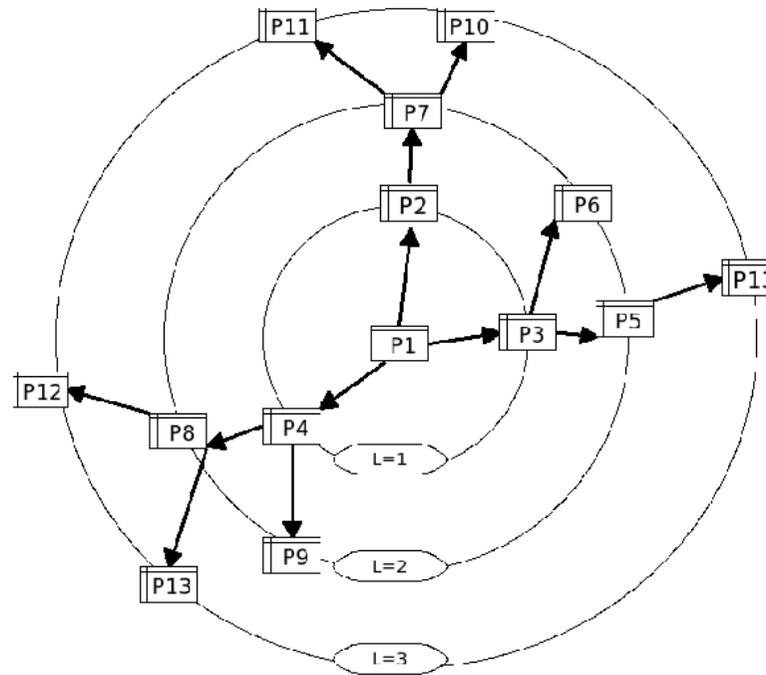


Single Trusted Page

- $t_1=1$; $t_i = 0$ $i \neq 1$
- Always jump to p_1 when bored
- We can place a **lowerbound** on being within L of P_1
- $N_L(p_1)$ =pages reachable from p_1 in L links



Single Trusted Page





Lower bound guarantee: Single Trusted

- Theorem 1:
 - Assuming the trust vector $T^{(1)}$, the sum of the PageRank values of all L -neighbors of p_1 is at least d^{L+1} close to 1.. That is:

$$\sum_{p_i \in N_L(p_1)} r_i \geq 1 - d^{L+1}$$



Lower bound guarantee: General Case

- The RankMass of the L-neighbors of the group of all trusted pages G , $N_L(G)$, is at least d^{L+1} close to 1.
That is:

$$\sum_{p_i \in N_L(G)} r_i \geq 1 - d^{L+1}$$



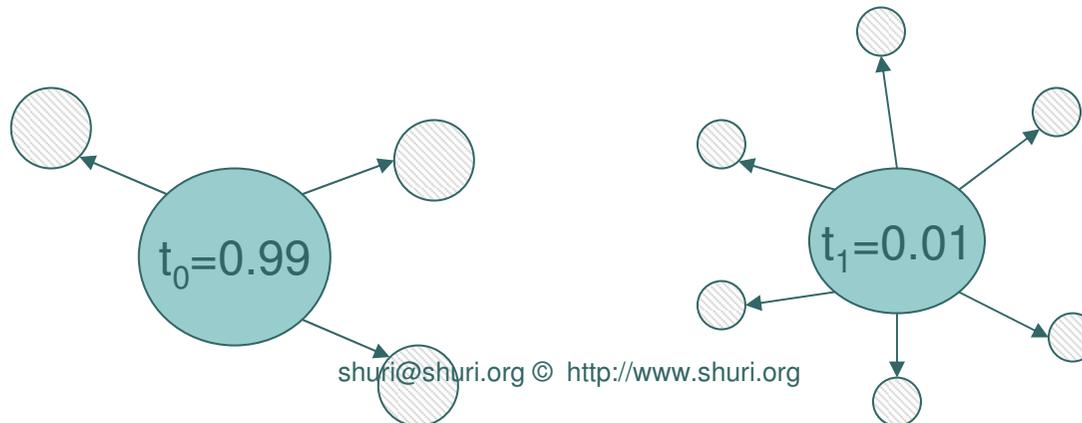
The L-Neighbor Crawler

1. $L := 0$
2. $N[0] = \{p_i | t_i > 0\}$ // Start with the trusted pages
3. While ($\square < d^{L+1}$)
 1. Download all uncrawled pages in $N[L]$
 2. $N[L + 1] = \{\text{all pages linked to by a page in } N[L]\}$
 3. $L = L + 1$



But what about efficiency?

- L-Neighbor similar to BFS
- L-Neighbor **simple and efficient**
- May wish to **prioritize further** certain neighborhoods first
- **Page level prioritization.**



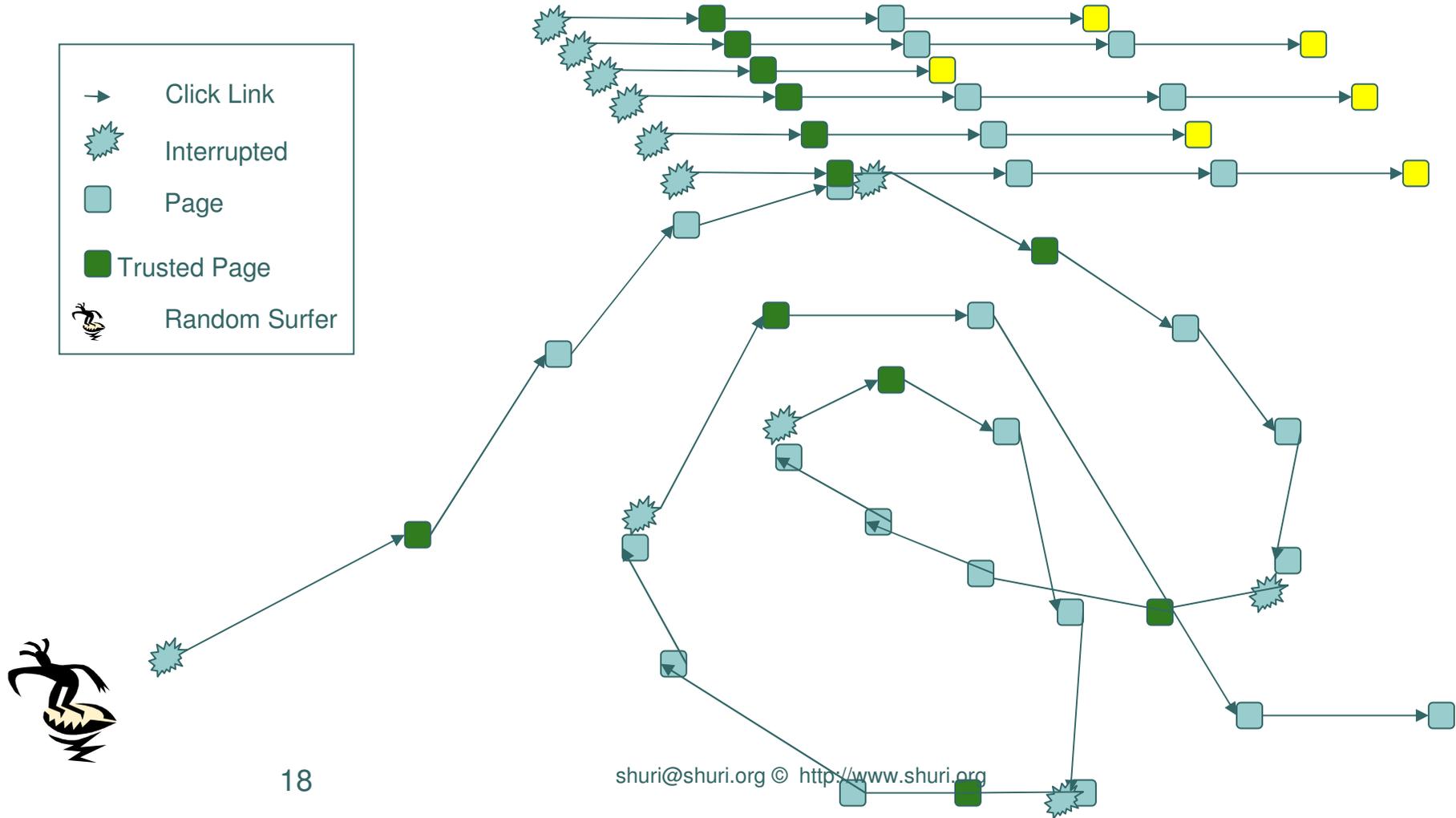


Page Level Prioritizing

- We want a more **fine-grained page-level** priority
- The idea:
 - Estimate PageRank on a **page basis**
 - High priority for pages with a high estimate of PageRank
- We **cannot calculate exact PageRank**
- **Calculate PageRank lower bound of undownloaded pages**



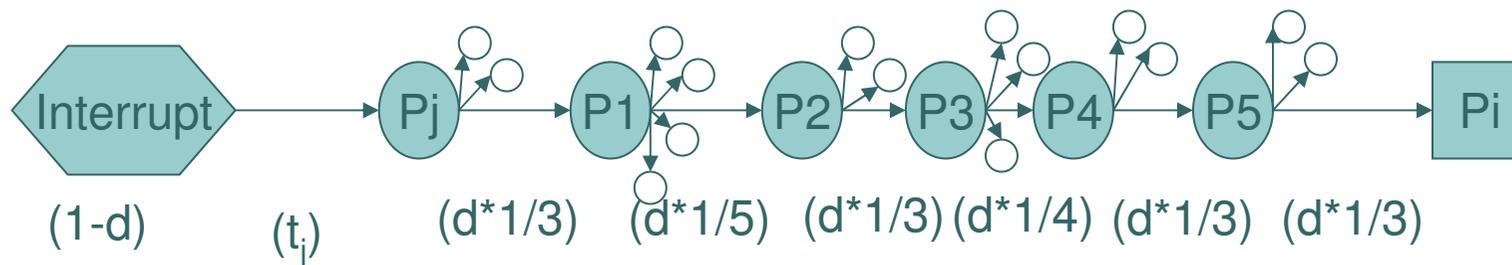
Probability of being at Page P





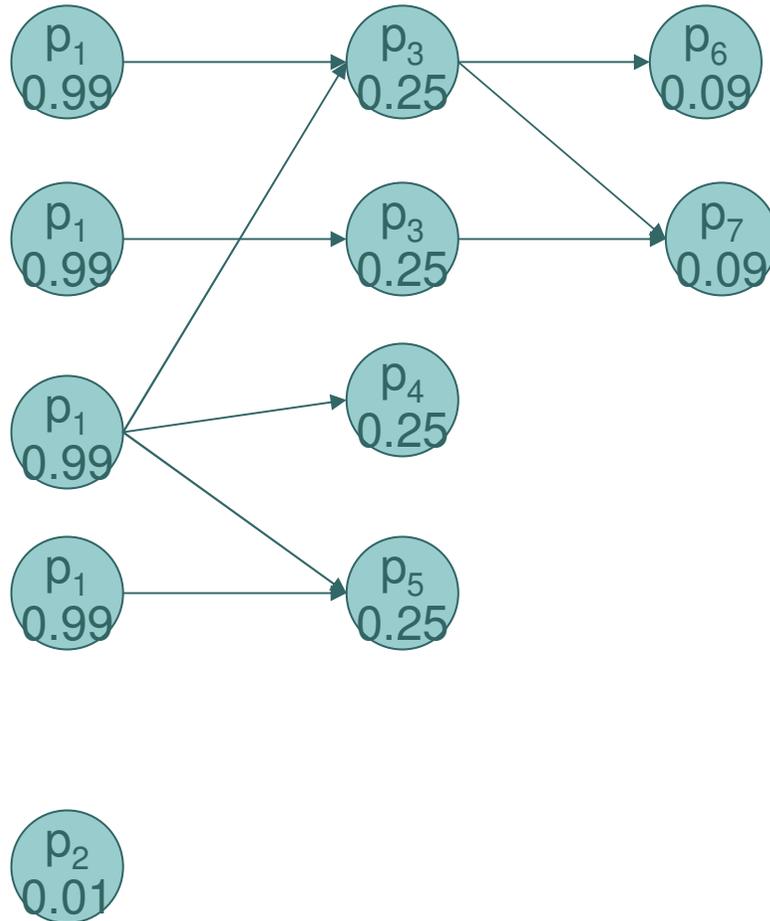
Calculating PageRank Lower Bound

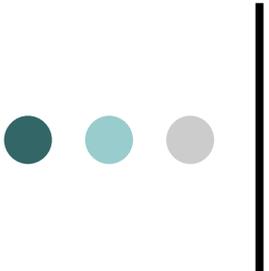
- PageRank(p) = Probability Random Surfer in p
- Breakdown path by "interrupts", jumps to a trusted page
- Sum up all paths that start with an interrupt and end with p





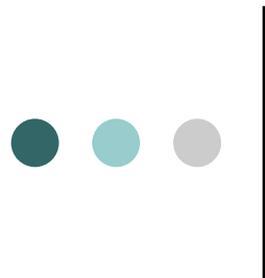
RankMass Basic Idea





RankMass Crawler: High Level

- But that sounds complicated?!
- Luckily we don't need all that
- Based on this idea:
 - Dynamically **update lower bound** on PageRank
 - **Update total RankMass**
 - **Download** page with **highest** lower bound



RankMass Crawler (Shorter)

- Variables:
 - CRM: RankMass lower bound of crawled pages
 - rm_i : Lower bound of PageRank of p_i .
- RankMassCrawl()
 - CRM = 0
 - $rm_i = (1 - d)t_i$ for each $t_i > 0$
 - While (CRM < $1 - \epsilon$):
 - Pick p_i with the largest rm_i .
 - Download p_i if not downloaded yet
 - CRM = CRM + rm_i
 - Foreach p_j linked to by p_i :
 - $rm_j = rm_j + d/c_i rm_i$
 - $rm_i = 0$



Experimental Setup

- HTML files only
- Algorithms simulated over web graph
- Crawled between Dec' 2003 and Jan' 2004
- 141 million URLs span over 6.9 million host names
- 233 top level domains.

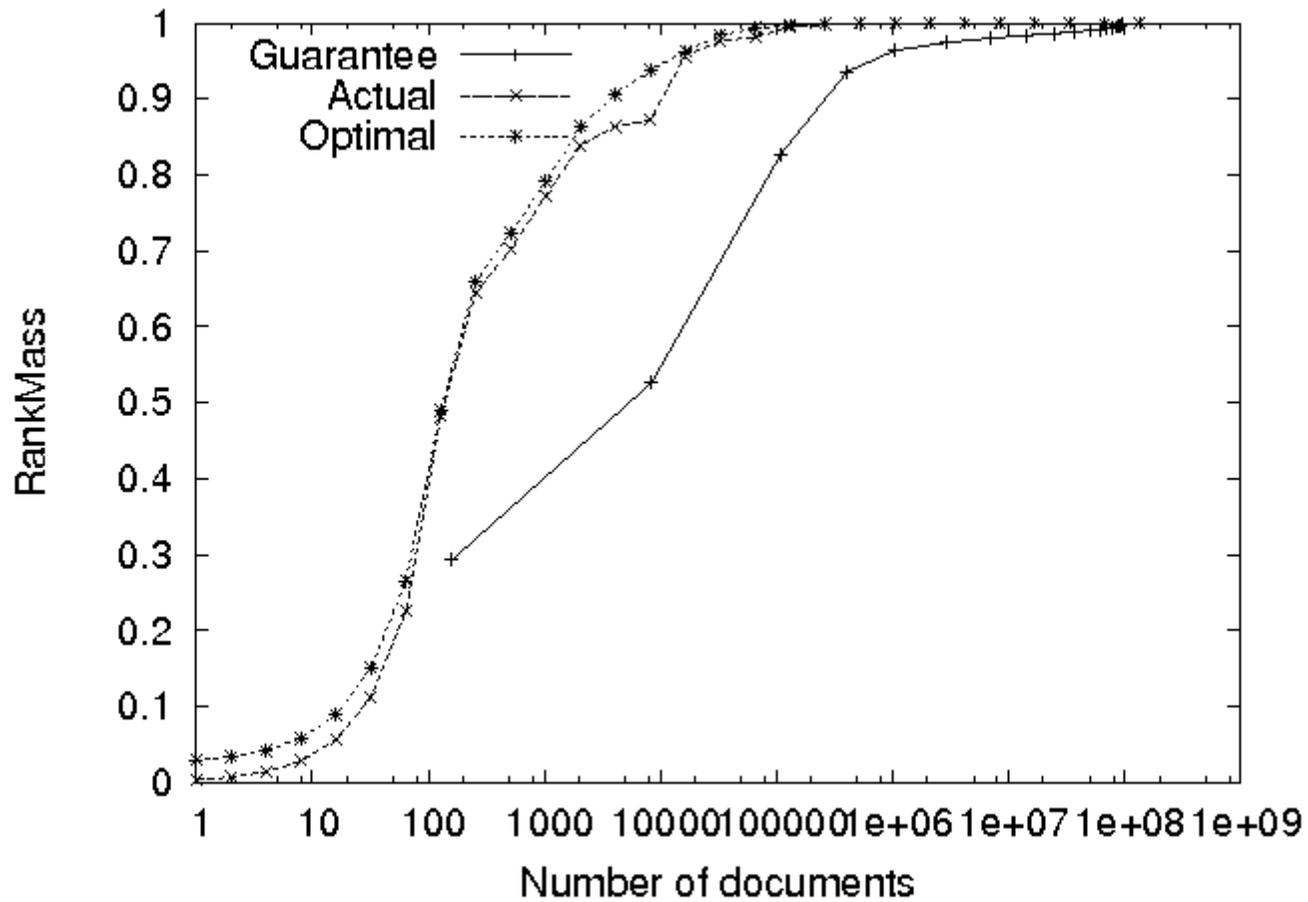


Metrics Of Evaluation

1. How much RankMass is collected during the crawl
2. How much RankMass is "known" to have been collected during the crawl
3. How much computational and performance overhead the algorithm introduces.

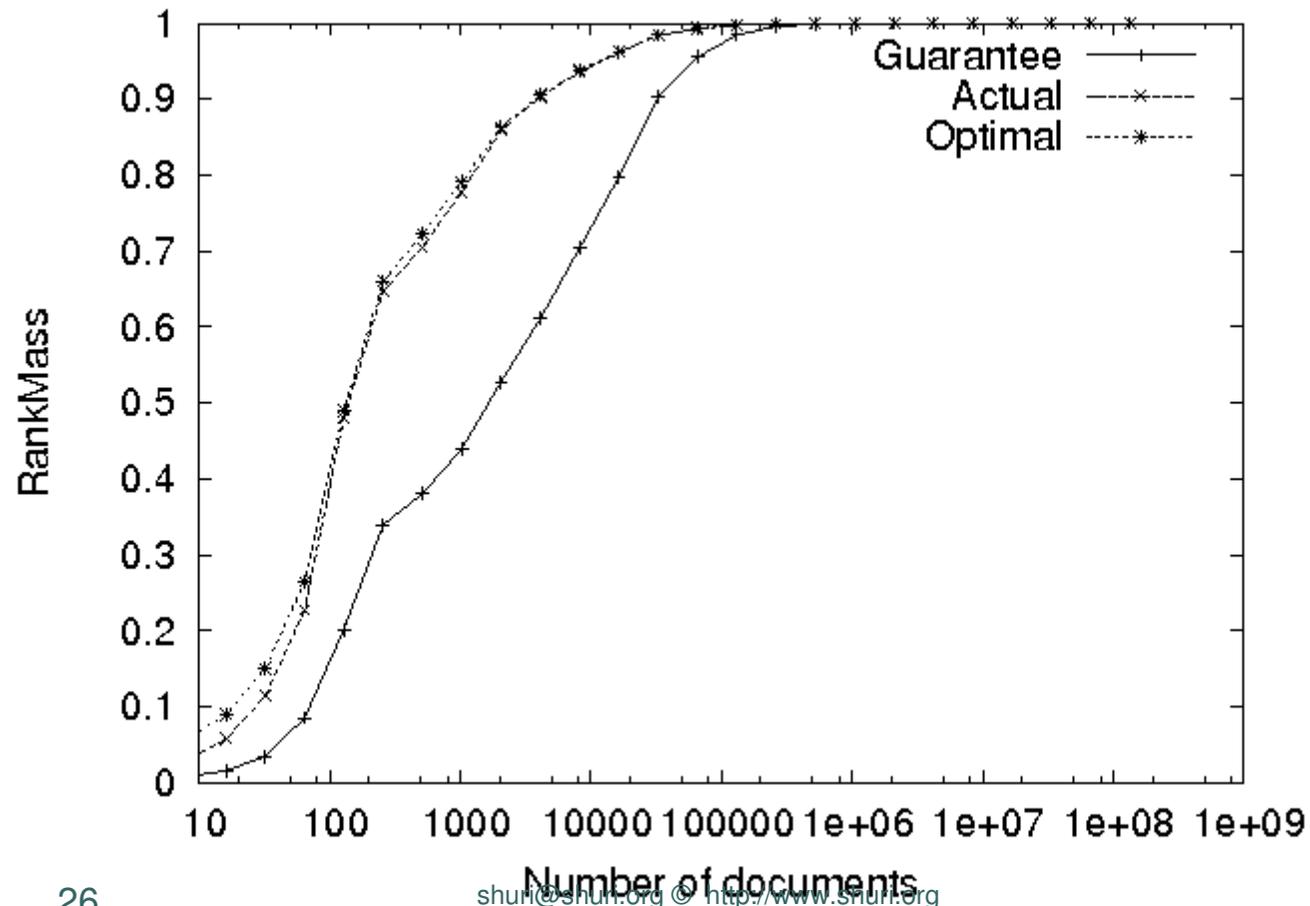


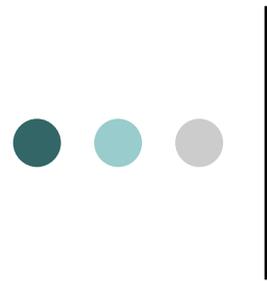
L-Neighbor





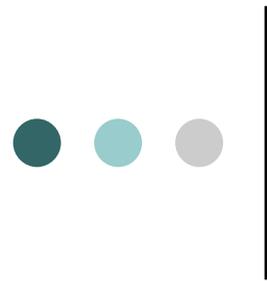
RankMass





Algorithm Efficiency

Algorithm	Downloads required for above 0.98% guaranteed RankMass	Downloads required for above 0.98% actual RankMass
L-Neighbor	7 million	65,000
RankMass	131,072	27,939
Windowed-RankMass	217,918	30,826
Optimal	27,101	27,101



Algorithm Running Time

Window	Hours	Number of Iterations	Number of Documents
L-Neighbor	1:27	13	83,638,834
20%- Windowed	4:39	44	80,622,045
10%- Windowed	10:27	85	80,291,078
5%- Windowed	17:52	167	80,139,289
RankMass	25:39	Not comparable	10,350,000



Thank you

- o Thank you





Greedy vs Simple

- L-Neighbor is **simple**
- RankMass is very **greedy**.
- Update **expensive**: random access to web graph
- Compromise?
- Batching
 - **downloads** together
 - **updates** together

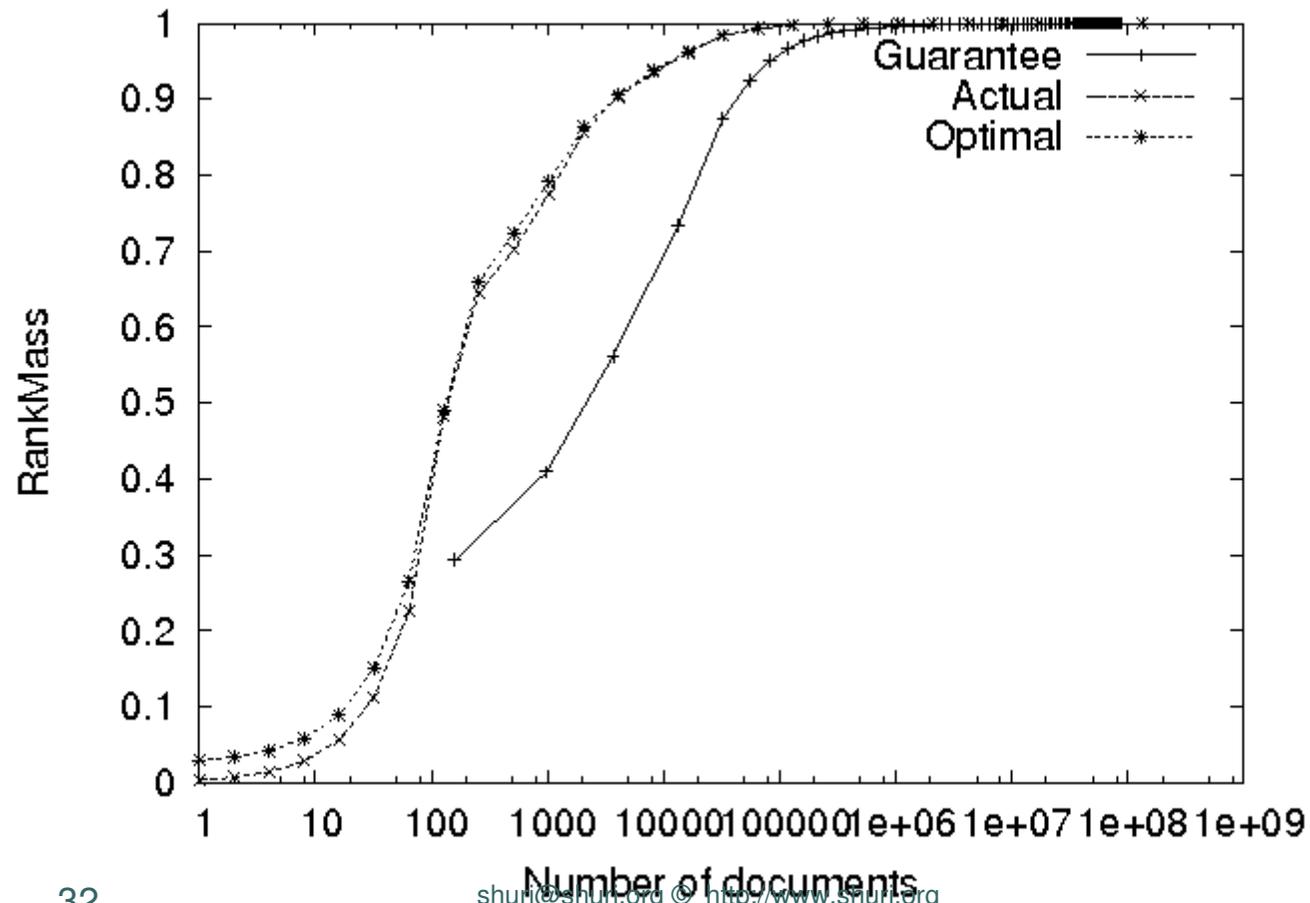


Windowed RankMass

- Variables:
 - CRM: RankMass lower bound of crawled pages
 - rm_i : Lower bound of PageRank of p_i .
- Crawl()
 - $rm_i = (1 - d)ti$ for each $ti > 0$
 - While (CRM < $1 - \epsilon$):
 - Download **top window%** pages according to rm_i
 - Foreach page $p_i \in D_c$
 - $CRM = CRM + rm_i$
 - Foreach p_j linked to by p_i :
 - $rm_j = rm_j + d/c_i rm_i$
 - $rm_i = 0$

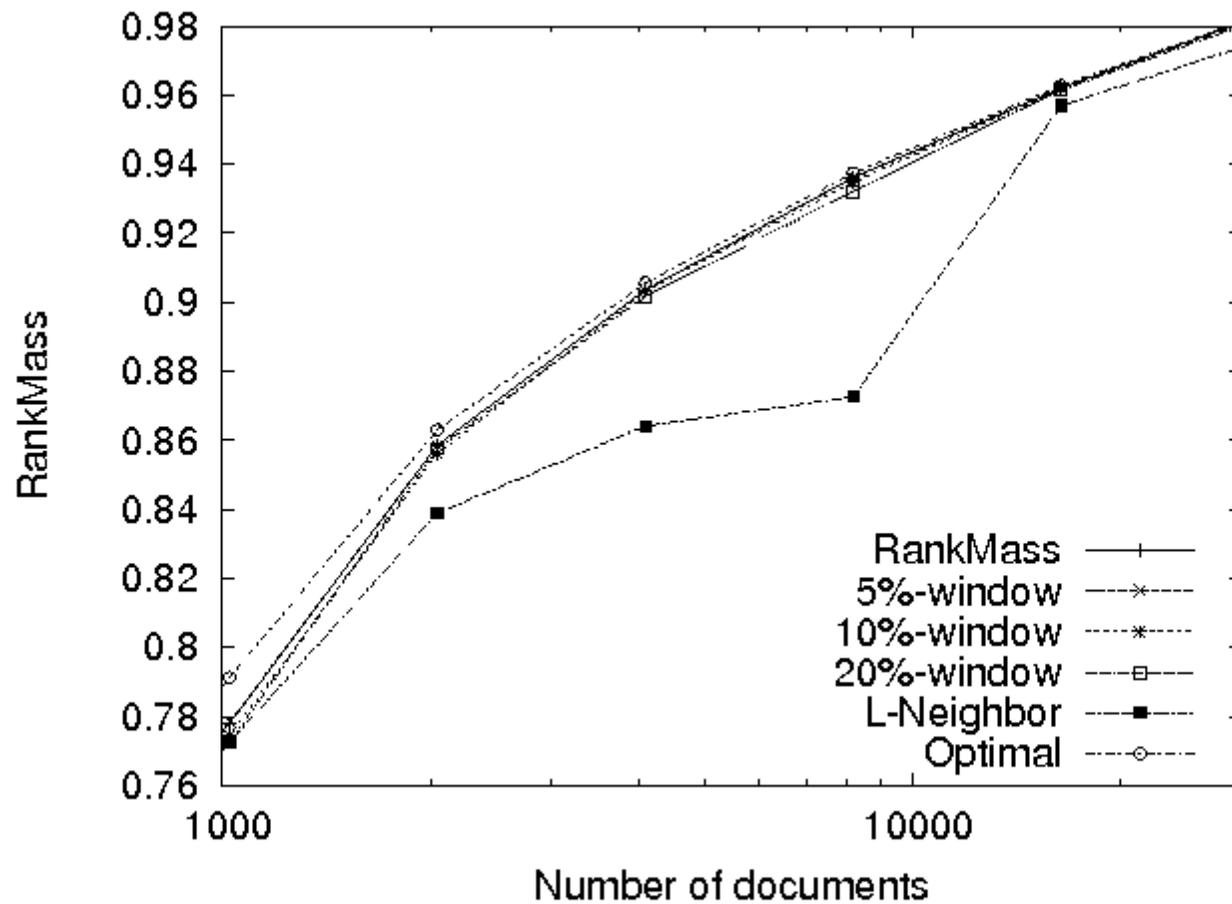


Windowed RankMass





Window Size





RankMass Lower Bound

- Lower bound given a single trusted page

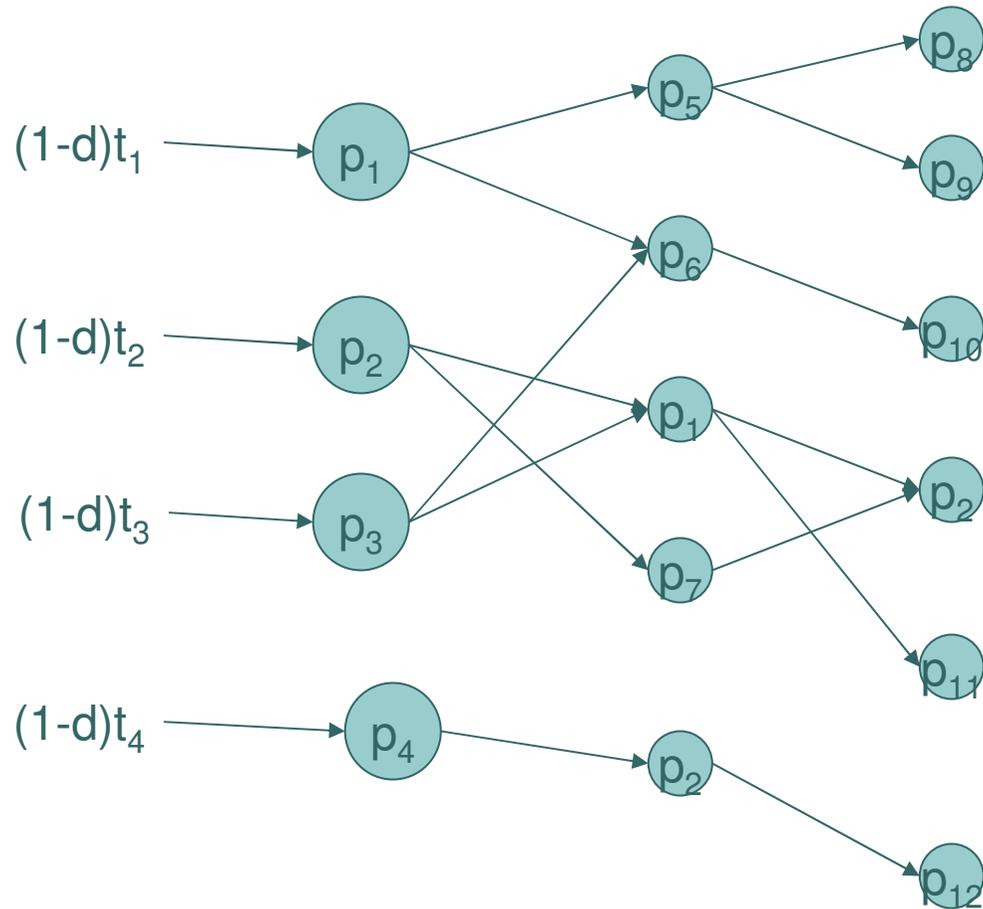
$$\sum_{p_i \in N_L(p_1)} r_i \geq 1 - d^{L+1}$$

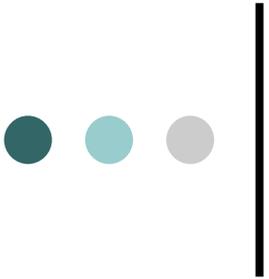
- That's the basis of the crawling algorithm with a coverage guarantee
- Extension: Given a set of trusted pages G

$$\sum_{p_i \in N_L(G)} r_i \geq 1 - d^{L+1}$$



Power Method of Calculating PageRank





RankMass Algorithm

- **Variables:**
 - **UnexploredPaths:** *List of unexplored paths and their path probabilities*
 - **sumPathProb_i:** *Sum of the probabilities of all unexplored paths leading to p_i*
 - **r_i:** *Partial sum of the probability of being in p_i*
- **RankMassCrawl()**
 - **// Initialize:**
 - **r_i = 0** for each i // *Set initial probability sum to be zero.*
 - **UnexploredPaths = {}** // *Start with empty set of paths.*
- **Foreach (t_i > 0):** // *Add initial paths of jumping to a trusted page and*
 - **Push [path: {p_i}, prob: (1 - d)t_i]** to **UnexploredPaths** // *the probability of the random jump.*
 - **sumPathProb_i = (1 - d)t_i** // *For every trust page p_i, we currently have only one path {p_i}*



RankMass Algorithm

- While ($\sum_i r_i < 1 - \epsilon$):
- Pick p_i with the largest sumPathProb_i . // Get the page with highest sumPathProb_i .
- Download p_i if not downloaded yet // Crawl the page.
- // Now expand all paths that end in p_i
- $\text{PathsToExpand} = \text{Pop all paths ending with } p_i$ // Get all the paths leading to p_i , from UnexploredPaths
- Foreach p_j linked to from p_i // and expand them by adding p_i 's children to the paths.
 - Foreach $[\text{path}, \text{prob}] \in \text{PathsToExpand}$
 - $\text{path}' = \text{path} \cdot p_j$ // Add the child p_j to the path,
 - $\text{prob}' = d/c_{ij} \cdot \text{prob}$ // compute the probability of this expanded path,
 - Push $[\text{path}', \text{prob}']$ to UnexploredPaths // and add the expanded path to UnexploredPaths .
 - $\text{sumPathProb}_j = \text{sumPathProb}_j + d/c_{ij} \cdot \text{sumPathProb}_i$ // Add the path probabilities of the newly added paths to p_j .
- Add the probabilities of just explored paths to r_i
- $r_i = r_i + \text{sumPathProb}_i$ // We just explored all paths to p_i . Add their probabilities
- $\text{sumPathProb}_i = 0$ // to r_i .