

Making Eigenvector-based Reputation Systems Robust to Collusion

**Hui Zhang¹, Ashish Goel², Ramesh Govindan¹, Kahn Mason²,
Benjamin Van Roy²**

¹University of Southern California

²Stanford University

Outline

- Research motivation.
- PageRank algorithm : a brief introduction.
- Study of PageRank's robustness to collusion.
- Adaptive-resetting: make PageRank robust to collusion.
- Conclusion & future works.

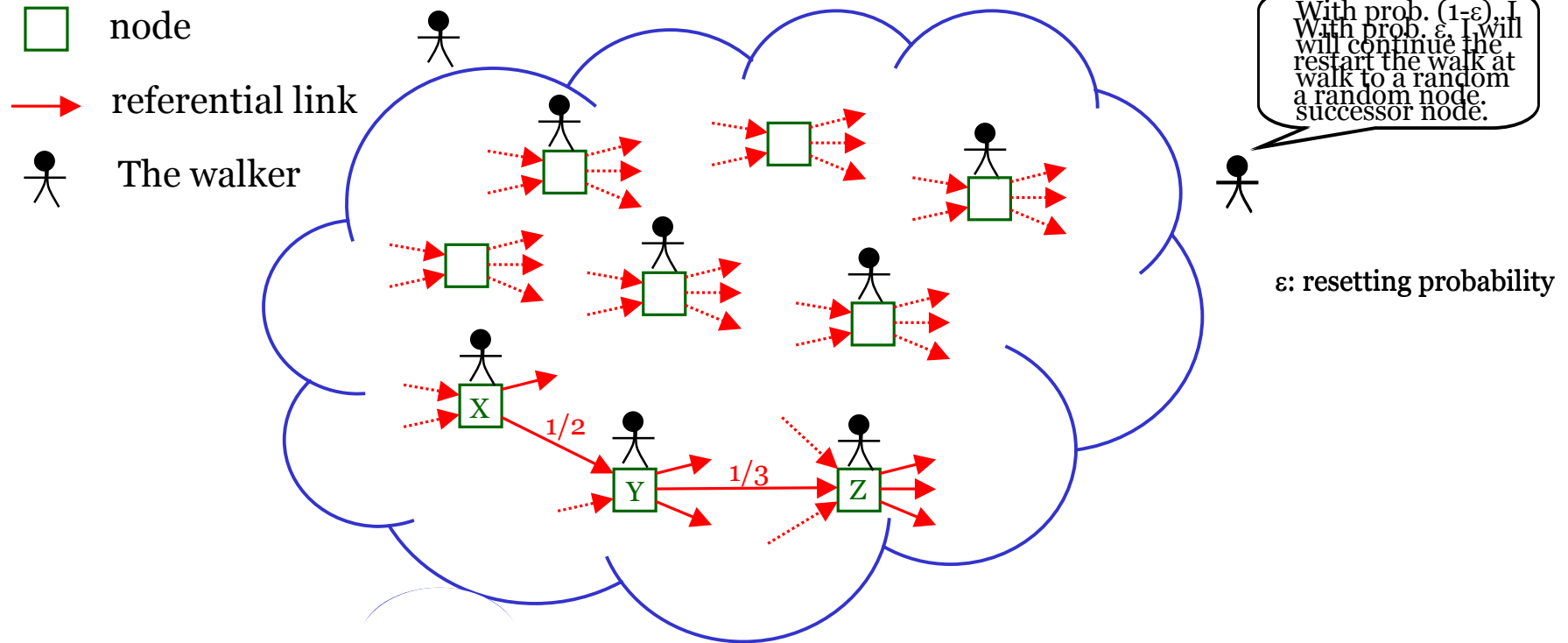
Research motivation

- Build reputation in large-scale systems
 - ❑ P2P file sharing systems
 - ❑ Blogging communities
 - ❑ Networked gaming, ..., etc.
- Collusion-proofness is an essential criterion in evaluating a rating scheme.

PageRank [Brin1998]

- An *eigenvector-based* rating scheme to rank hypertext documents on the WWW.
- An iterative algorithm to calculate the importance of a web page based on the importance of its parent pages.
- Can be applied to other systems than WWW.

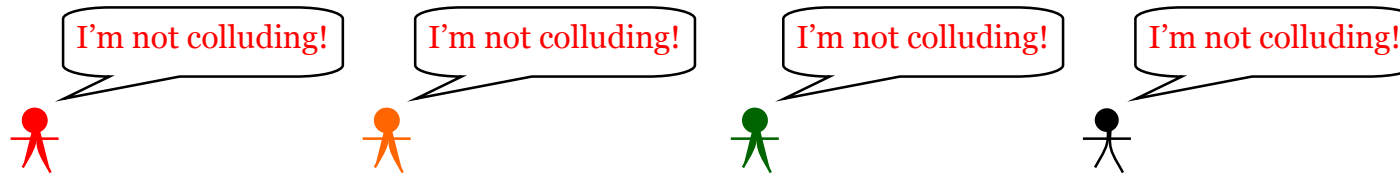
PageRank: random walk model



- As time goes on, the expected percentage of steps the walker is at each node v converges to the PageRank weight $PR(v)$.

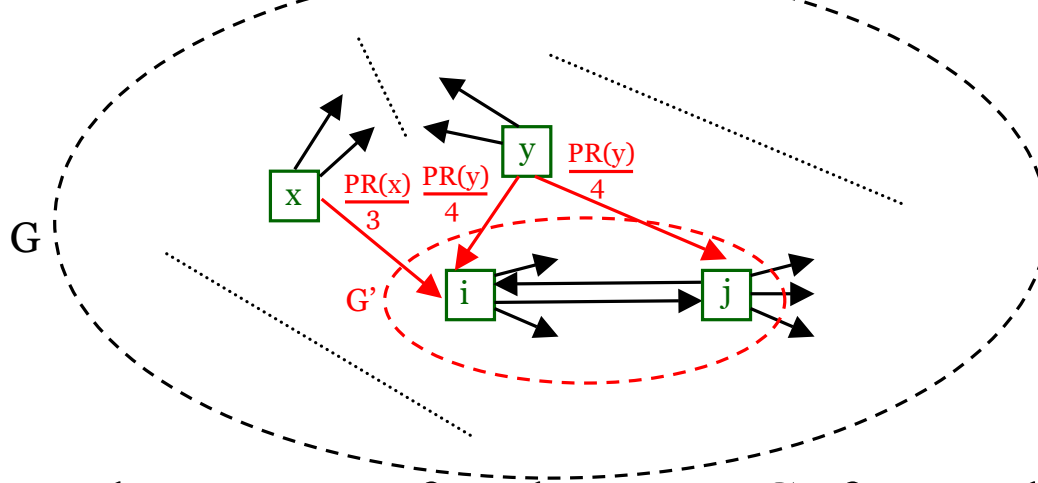
PageRank: is it collusion-proof?

- Can a node easily boost its rank by manipulating its outgoing links with others'?



Amp(G): a metric on group collusion

ϵ : resetting probability



$$W_G(G') = PR(i) + PR(j)$$

← real group weight

$$W_{in}(G') = \frac{PR(x)}{3} (1-\epsilon) + \frac{PR(y)}{2} (1-\epsilon) + \frac{2}{N} (1-W(G'))\epsilon$$

← "actual" group weight

- In the system of node group G , for a subgroup G' ,

the amplification factor $Amp(G') = \frac{W_G(G')}{W_{in}(G')}$

- $$W_G(G') = \sum_{i: i \in G'} PR(i)$$

- $$W_{in}(G') = \sum_{(i,j): i \notin G', j \in G', \exists i \rightarrow j} \frac{PR(i)}{out(i)} (1-\epsilon) + \frac{|G'|}{|G|} (1-W_G(G'))\epsilon$$

Theorem on Amp

- In the original PageRank system,

$$\forall G' \subseteq G, \text{Amp}(G') \leq \frac{2}{\varepsilon}$$

where ε is the resetting probability.

Specifically, when $\frac{|G'|}{|G|} \ll 1$, $\text{Amp}(G') \leq \frac{1}{\varepsilon}$

Two experimental topologies

- W , a Web link topology
 - Contains the link structure of upwards of 80 million URLs.
 - Source: the Stanford WebBase.
- B , a weblog blogrolling topology
 - Contains the blogrolling structure of upwards of 72,000 blogs.
 - Source: *www.blogstreet.com* the XML -RPC weblog service.

Experiment 1: Collusion200

- Model a small number of web pages *simultaneously* colluding.
- Methodology:
 - 100 colluding groups of 200 nodes;
 - Each colluding group has the circle topology consisting of two nodes with adjacent ranks;
 - Arbitrarily chose node pairs originally ranked around 1000th, 2000th, ..., 100000th.
 - $\varepsilon = 0.15$.

Experiment result of *Collusion200* (I)

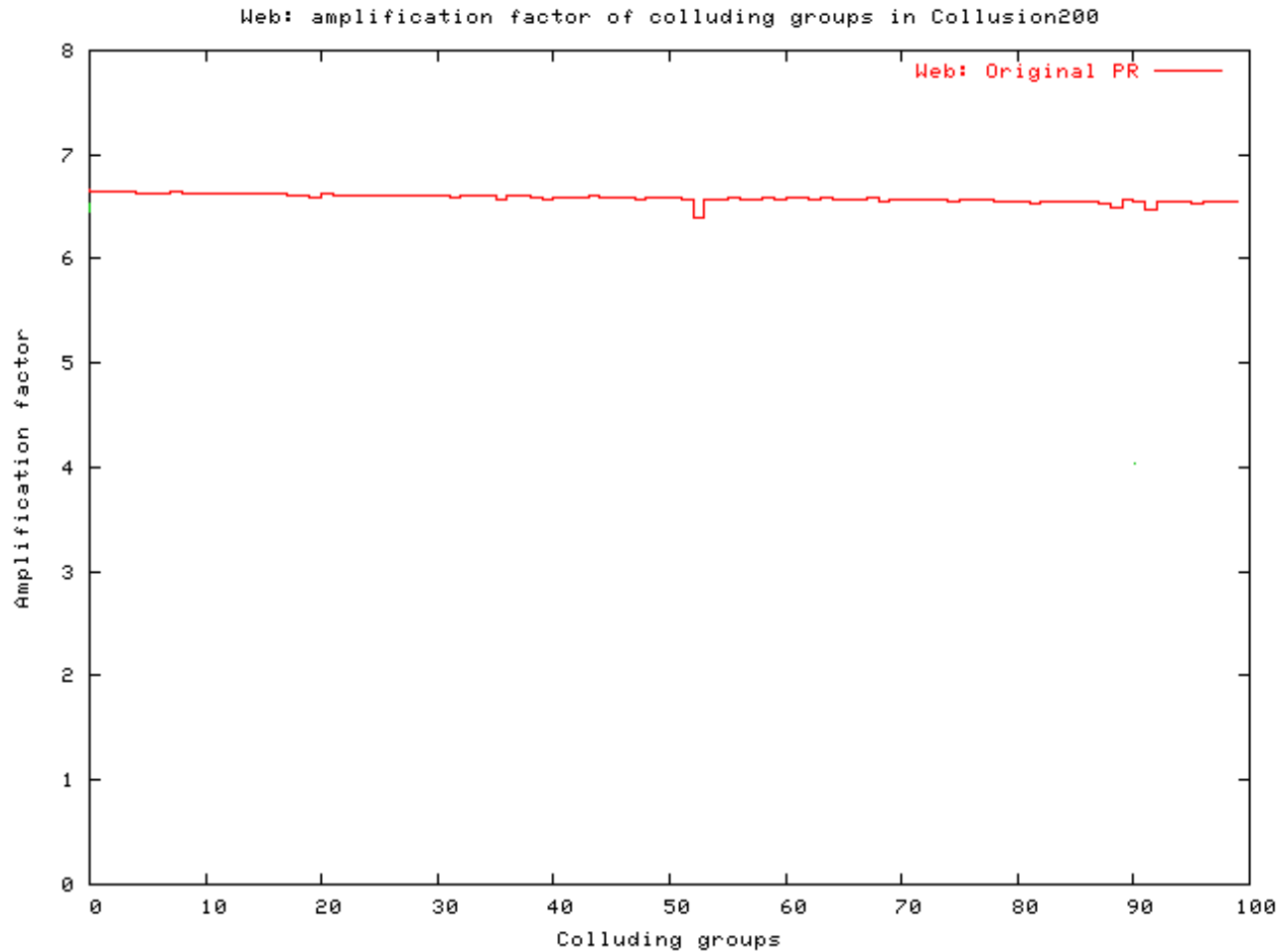


Figure 1: W - Amplification factors of the 100 colluding groups in *Collusion200*.

Experiment result of *Collusion200* (III)

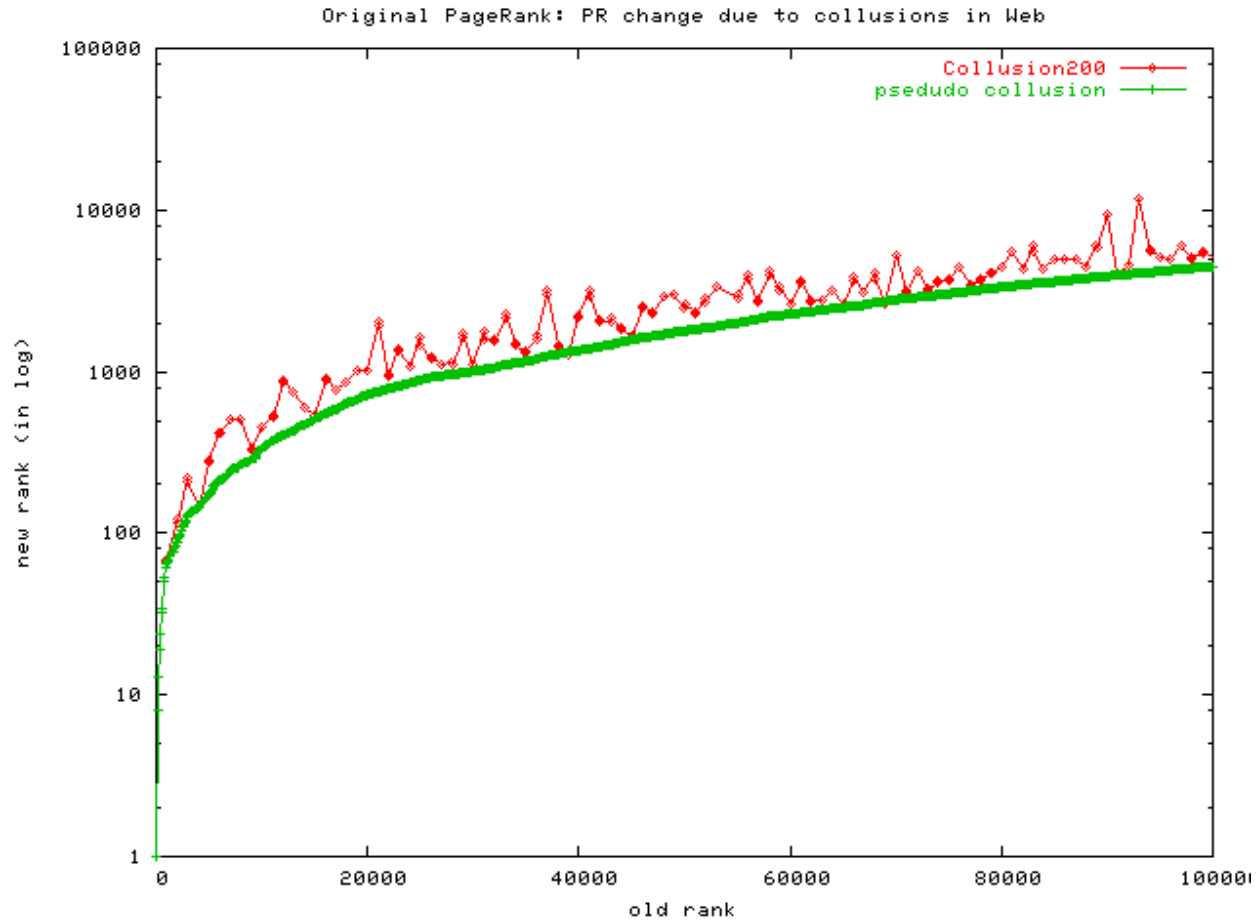


Figure 2: W – new PR rank after *Collusion200*.

There is a long flat portion...

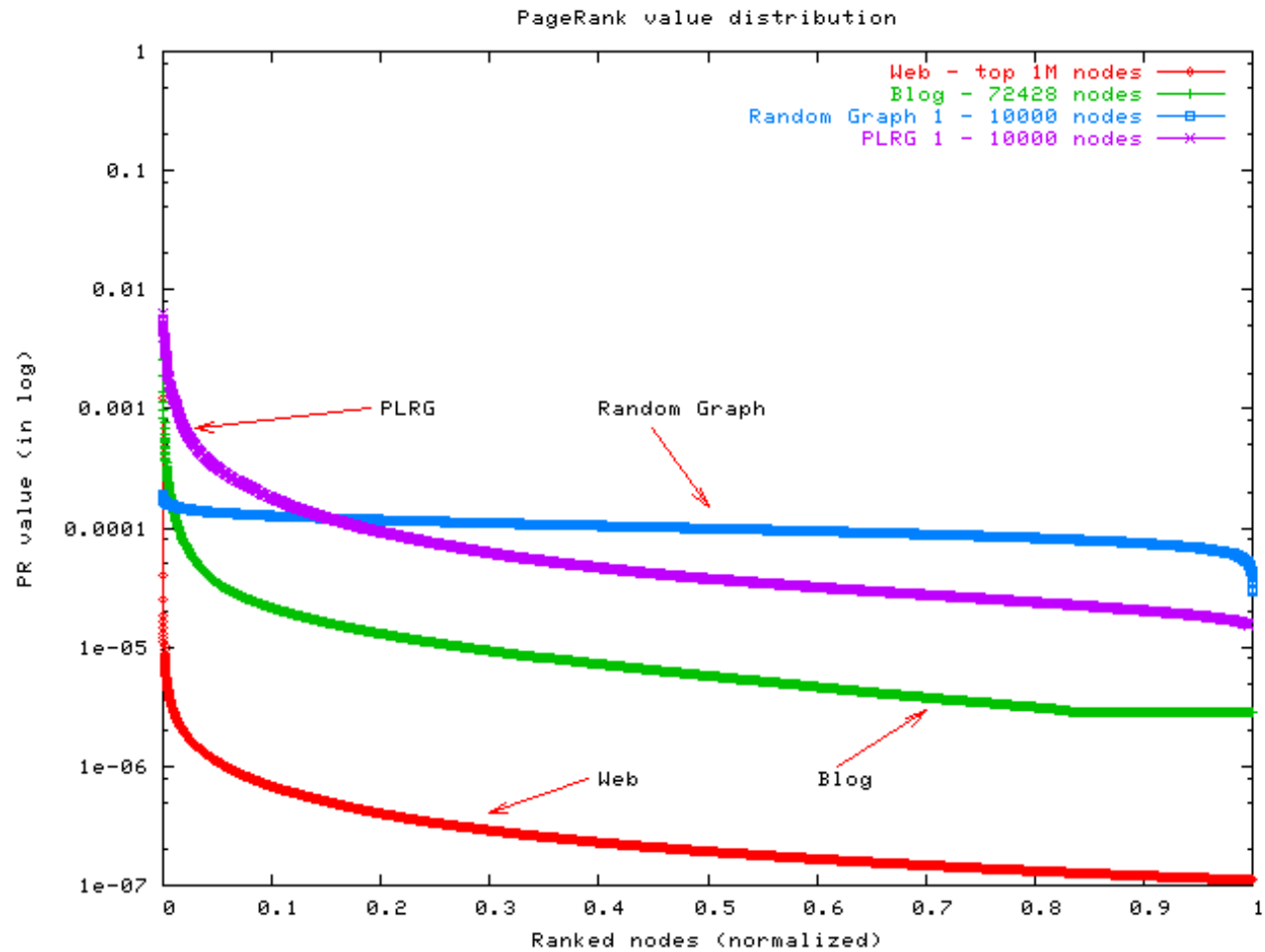


Figure 3: The PR weight distribution of 4 topologies.

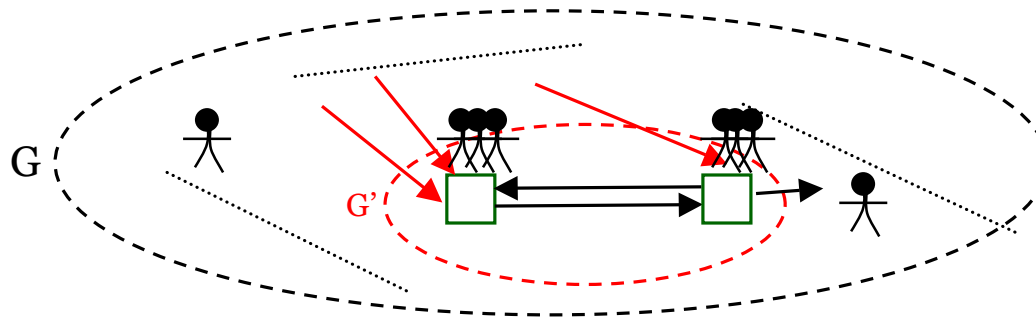
Next step: how to detect collusions?

- Theorem on Hardness.

Max $_{G' \subseteq G} \text{Amp}(G')$ is a NP-Hard problem.

An observation on collusion behaviors

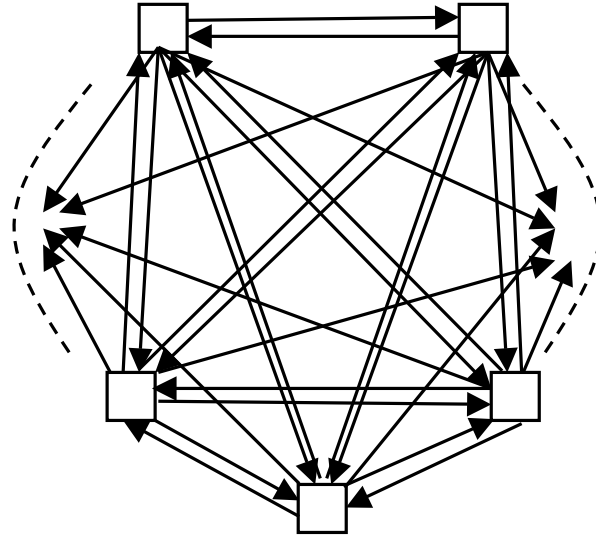
- To increase their PR weight, i.e., the stationary weight in the random walk, the colluding nodes will stall the random walk.



- When the resetting probability ε increases, the colluding nodes must suffer a significant drop in PR weight.
- Therefore, we expect the PR weight of colluding nodes to be highly correlated with $1/\varepsilon$ (the average walk length), while that of non-colluding nodes is relatively insensitive to the change in ε .

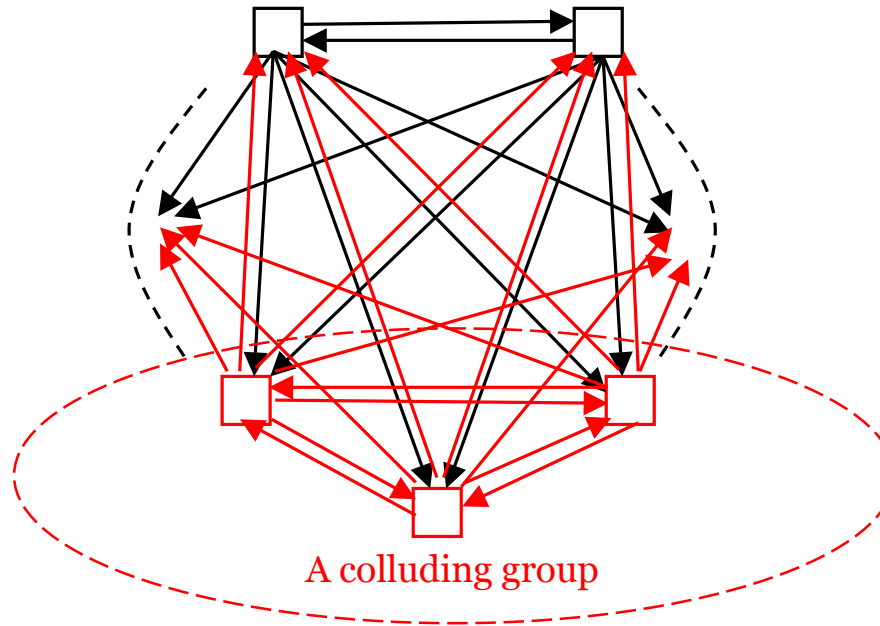
An intuitive example

- node
- referential link

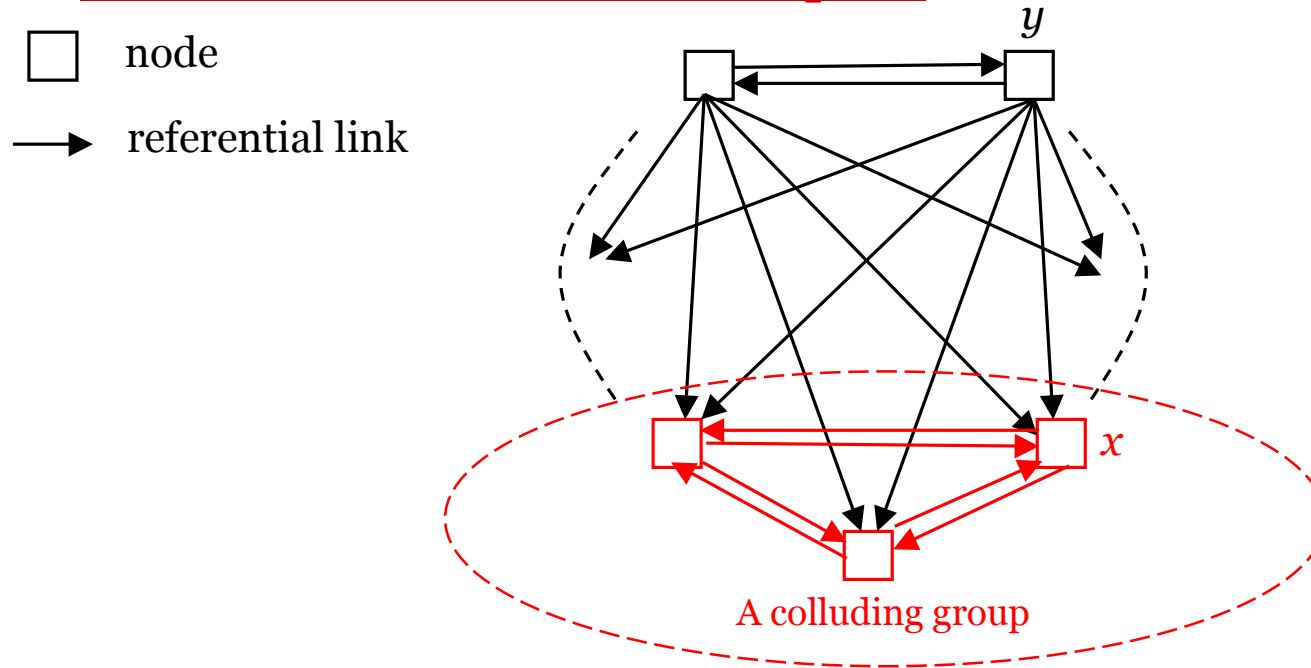


An intuitive example

- node
- referential link



An intuitive example



- A colluding node x : $\text{PR}(x) = \frac{1}{K + (N - K)\epsilon} \approx \frac{1}{N\epsilon}$, and $\text{co-co}(\text{PR}(x), 1/\epsilon) \approx 1$. (**co-co**: correlation coefficient)
- A non-colluding node y : $\text{PR}(y) = \frac{\epsilon}{K + (N - K)\epsilon} \approx \frac{1}{N}$, and $\text{co-co}(\text{PR}(y), 1/\epsilon) \approx 0$.

Adaptive-resetting scheme

- Part I – collusion detection:
 - Given the topology, calculate the PR vector under different ε values.
 - $\{\varepsilon\} = \{0.0375, 0.05, 0.075, 0.15, 0.3, 0.45, 0.6\}$, $\varepsilon_{default} = 0.15$.
 - Calculate the correlation coefficient between the curve of each node x 's PR weight and the curve of $1/\varepsilon$. Label it as $co-co(x)$.
- Part II – ε personalization:
 - Calculate each node x 's out-link personalized- $\varepsilon = F(\varepsilon_{default}, co-co(x))$.
 - Exponential function $F_{Exp} = \varepsilon_{default}^{(1.0-co-co(x))}$.
 - Linear function $F_{Linear} = \varepsilon_{default} + (0.5-\varepsilon_{default}) * co-co(x)$
 - The final PR weight vector is calculated with these personalized resetting values.

Experiment result of *Collusion200* (IV)

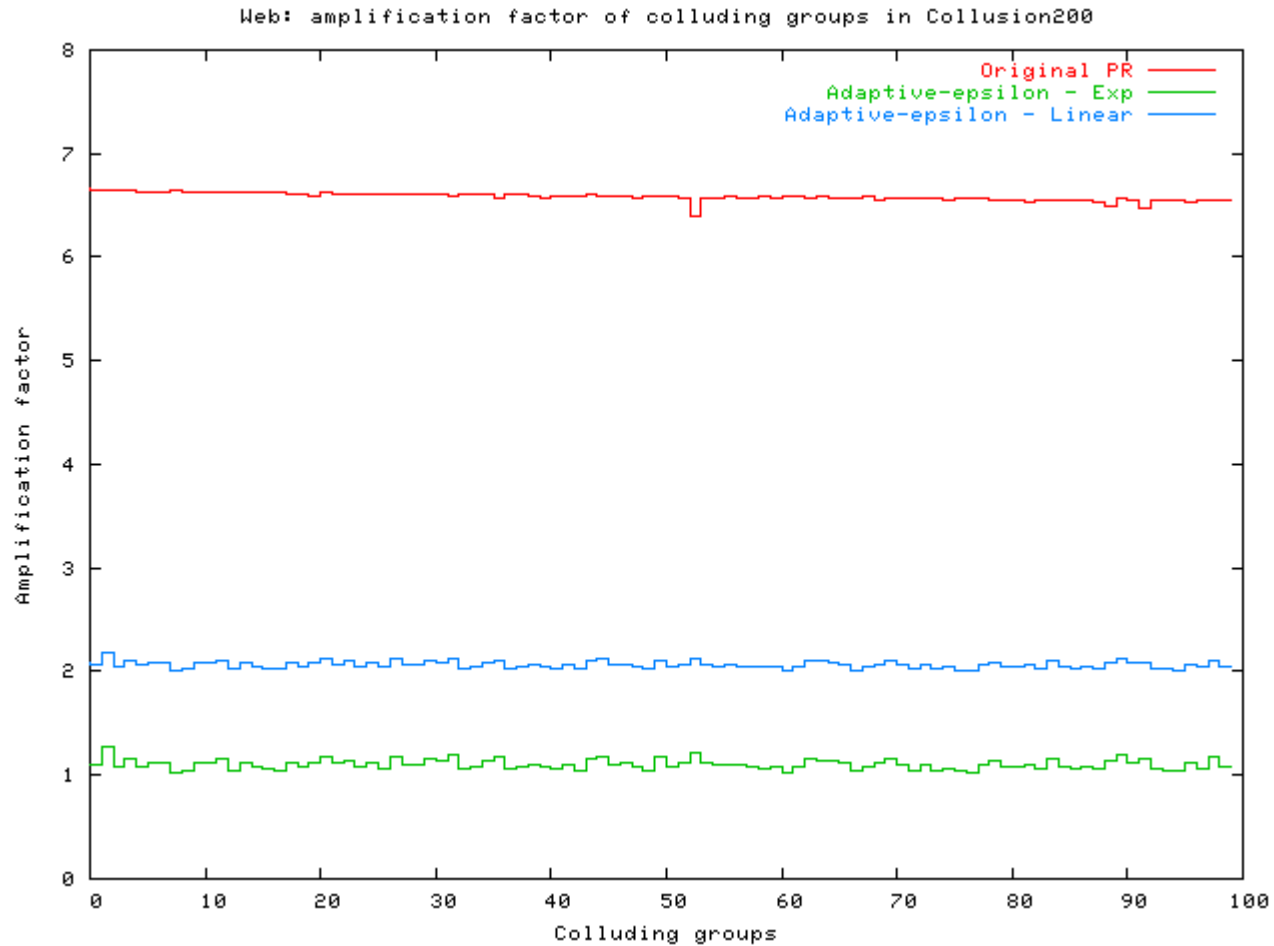


Figure 5: W - Amplification factors of the 100 colluding groups in *Collusion200*.

Experiment result of *Collusion200* (VI)

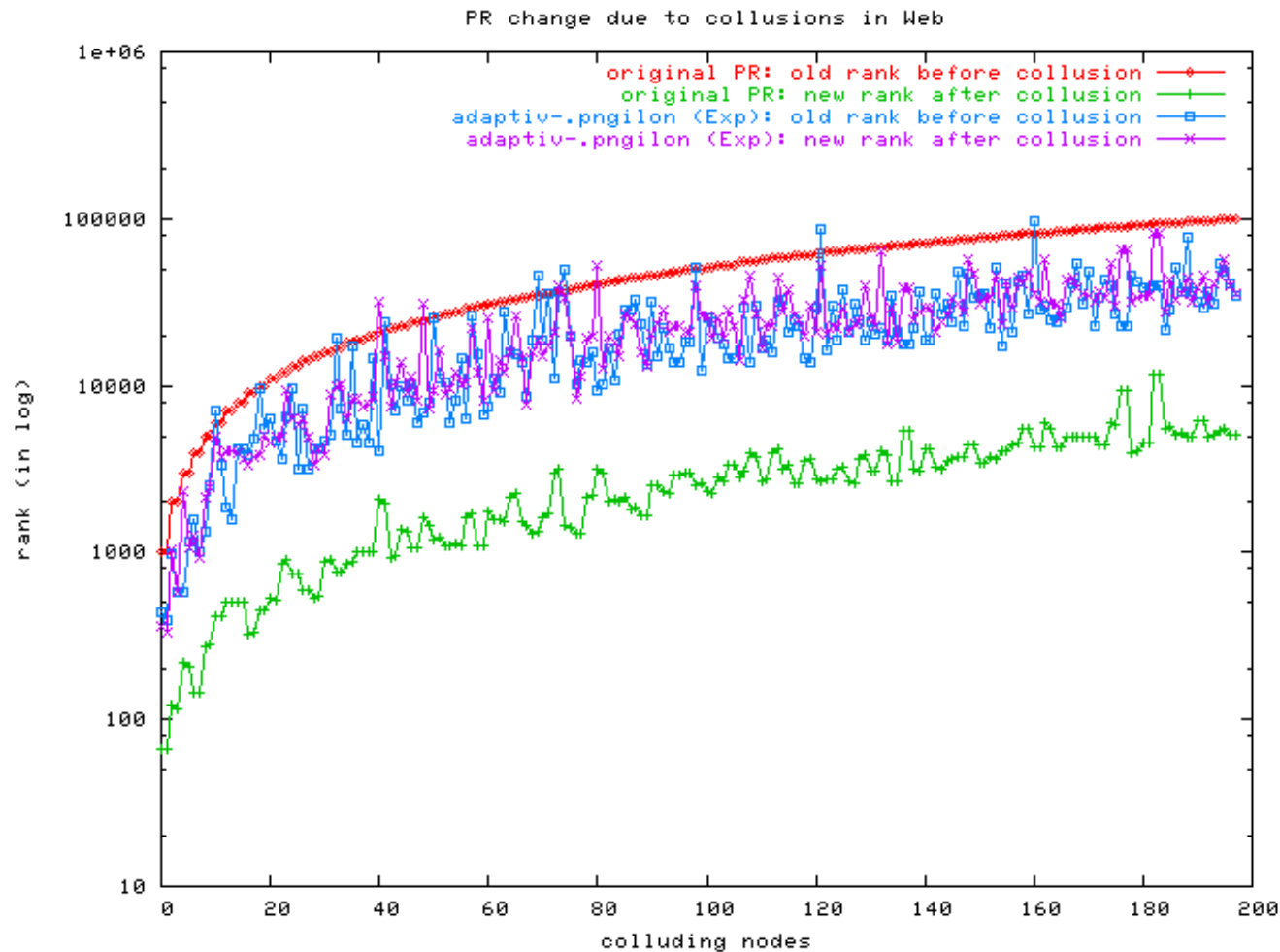
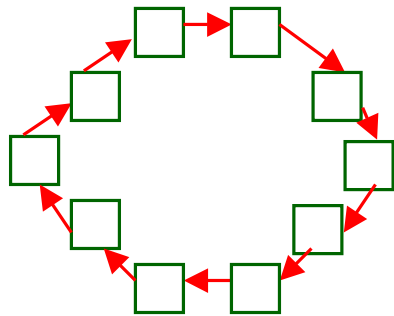


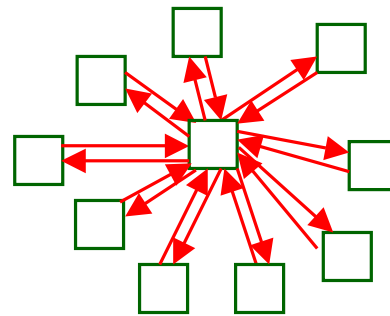
Figure 6: W – new PR rank after *Collusion200*.

Experiment 2: Collusion22

- Model various colluding subgraphs.
- Methodology:
 - 3 colluding groups:



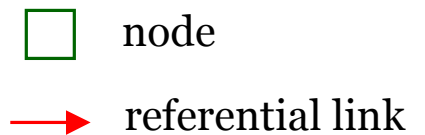
G1: 10-node ring



G2: 10-node star topology



G3: 2-node ring



Experiment result of *Collusion22* (I)

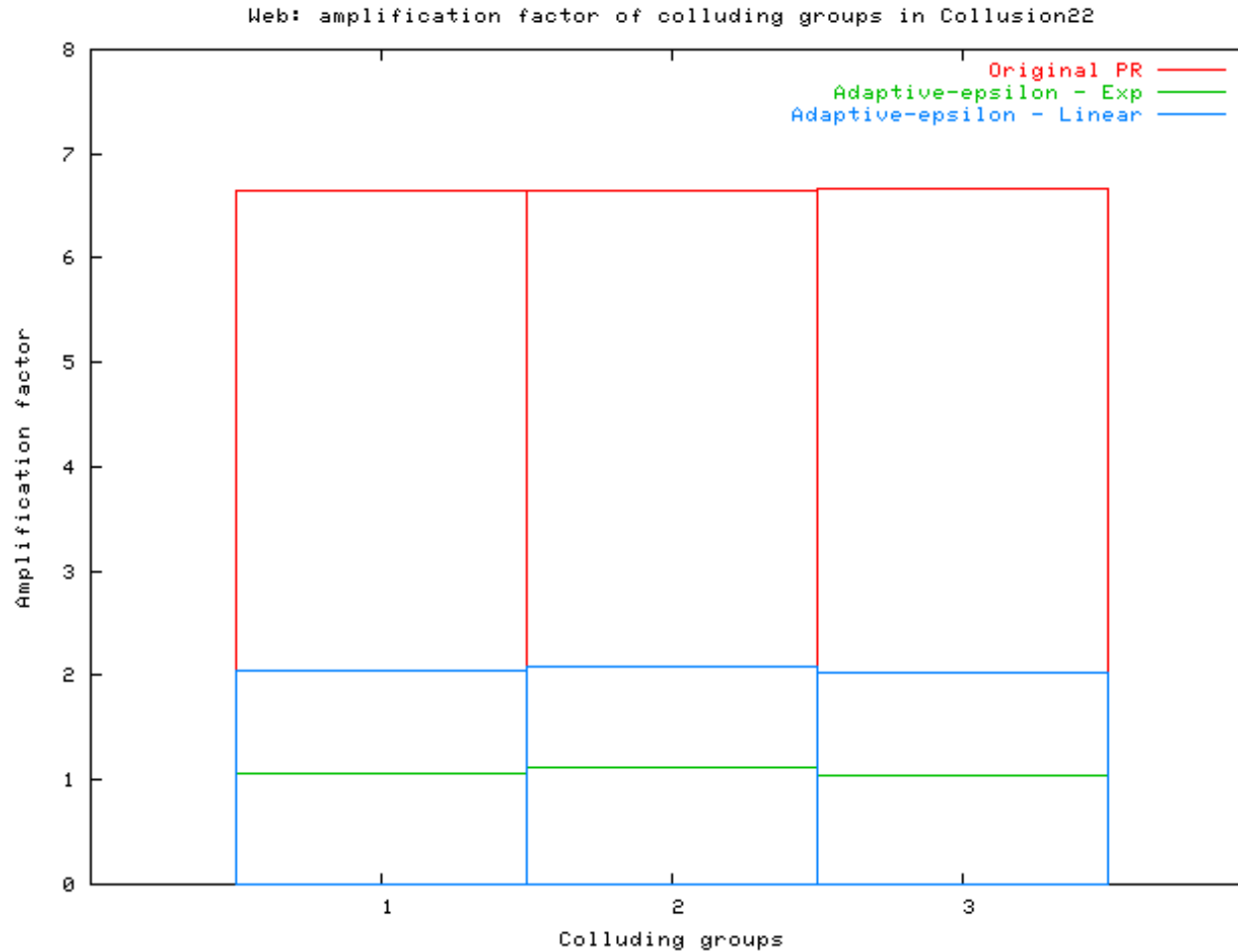


Figure 7: Amplification factors of the 3 colluding groups in *Collusion22*.

Experiment result of *Collusion22* (II)

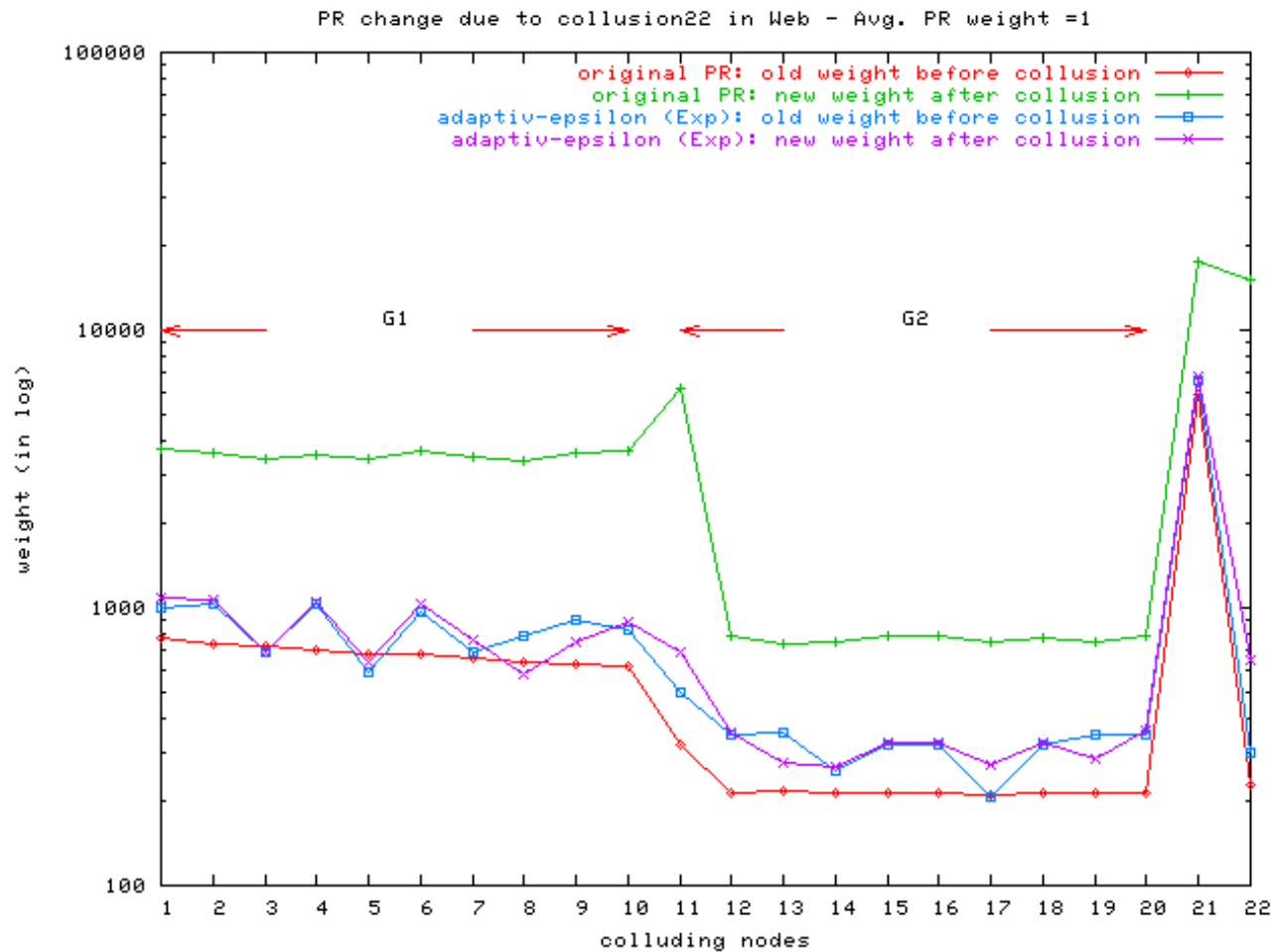


Figure 8: W – new PR weight after *Collusion22*.

New top-25 URL list in \mathcal{W}



Rank	Old list	New list
1	http://www.yahwa.com/	http://www.tucows.com/
2	http://messenger.yahoo.com/	http://www.yahoo.com/
3	http://www.tucows.com/	http://www.domaindirect.com/
4	http://www.domaindirect.com/	http://news.tucows.com/
5	http://news.tucows.com/	http://ispcentral.tucows.com/
6	http://ispcentral.tucows.com/	http://www.microsoft.com/
7	http://www.microsoft.com/	http://www.acme.com/software/thttpd
8	http://www.microsoft.com/info/copyright.htm	http://www.adobe.com/products/acrobat/readstep.html
9	http://www.adobe.com/products/acrobat/readstep.html	http://home.netscape.com/
10	http://home.netscape.com/	http://www.thecounter.com/
11	http://www.dun.com/	http://www.gendex.com/ged2html
12	http://www.worldwidemart.com/scripts	http://www.adobe.com/
13	http://www.acme.com/software/thttpd	http://www.worldwidemart.com/scripts
14	http://search.internet.com/	http://upload.tucows.com/contactus.html
15	http://upload.tucows.com/contactus.html	http://www.w3.org/
16	http://www.thecounter.com/	http://www.listbot.com/
17	http://www.listbot.com/	http://www.tucows.com/privacy.html
18	http://www.w3.org/	http://www.worldwidemart.com/scripts/faq/wwwboard
19	http://www.adobe.com/	http://www.microsoft.com/windows/w/default.htm
20	http://www.tucows.com/search.html	http://www.mga.gov/
21	http://www.tucows.com/privacy.html	http://www.hadi.com/
22	http://www.gendex.com/ged2html	http://www.rsa.org/
23	http://chl.levels.ac.uk/mikee/personal.html	http://search.internet.com/
24	http://www.achbar.com/misc/privacy.html	http://www.mca.gov/
25	http://www.achbar.com/homepage.html	http://chl.levels.ac.uk/mikee/personal.html

Table 1: The old and new top-25 list of \mathcal{W}

Conclusion & future works

- A collusion-proof rating scheme based on PageRank algorithm.
- Future works:
 - Formal analysis of the adaptive-resetting scheme.
 - Study of Web link structure evolution under PageRank within the framework of game theory.