

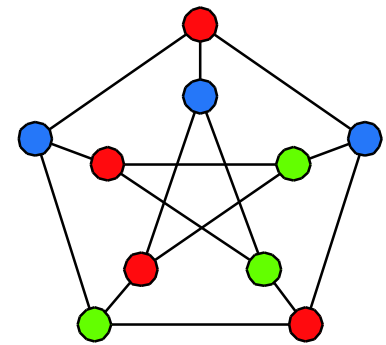
# Efficient PageRank Tracking in Evolving Networks

Naoto Ohsaka (UTokyo)

Takanori Maehara (Shizuoka University)

Ken-ichi Kawarabayashi (NII)

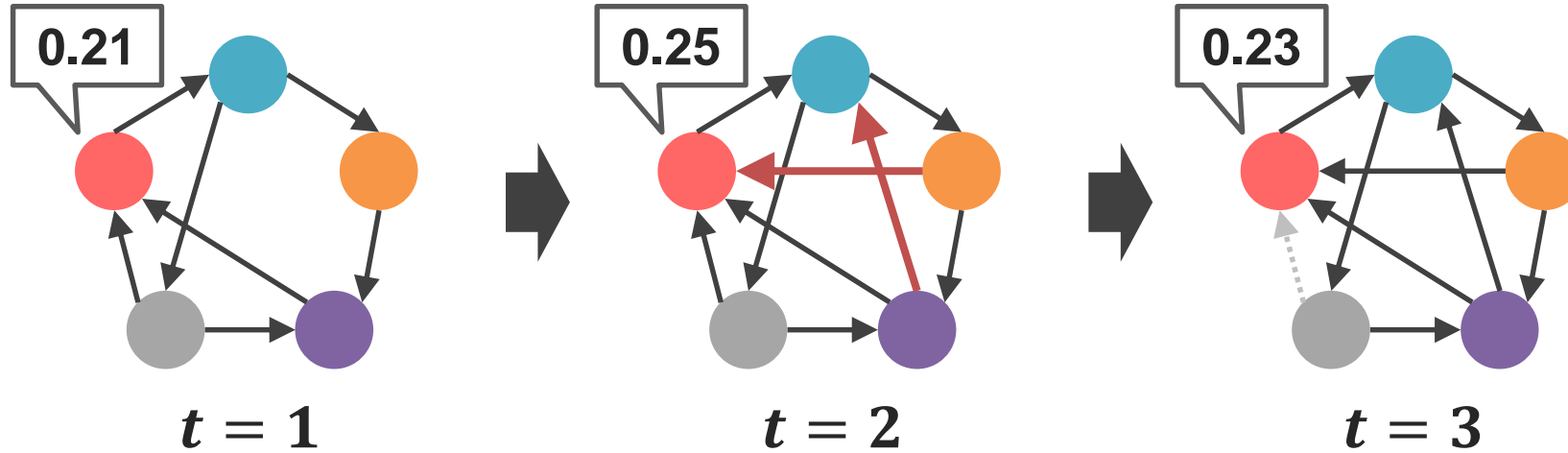
*ERATO Kawarabayashi Large Graph Project*



**ERATO**

## Background

# Personalized PageRank Tracking



## Applications

Search engine [Brin-Page. '98]

Spam detection

[Chung-Toyoda-Kitsuregawa. '09, '10]

User recommendation

[Gupta-Goel-Lin-Sharma-Wang-Zadeh. '13]

## Growth of real networks

	Size	Speed
WWW	60T	600K pages / s
Twitter	300M	5K tweets / s
Google+	700M	+19 users / s

# Existing Work for PageRank Tracking

	$m$ random edge insertions	<b>Scalability</b> Update time < 0.1s Error $\approx 10^{-9}$
Aggregation/Disaggregation <small>[Chien et al. '04]</small>	$\mathcal{O}(m S  \log 1/\epsilon)$	68M edges
Monte-Carlo <small>[Bahmani et al. '10]</small>	$\mathcal{O}(m + \log m / \epsilon^2)$	68M edges
Power method <small>naive method</small>	$\mathcal{O}(m^2 \log 1/\epsilon)$	11M edges

# Our Contribution

Propose a **simple**, **efficient**, & **accurate** method for Personalized PageRank tracking in evolving networks

	<b><math>m</math> random edge insertions</b>	<b>Scalability</b> Update time < 0.1s Error $\approx 10^{-9}$
<b>This work</b>	Ave. $\downarrow$ Max. out-deg <b><math>\mathcal{O}(m + \Delta \log m / \epsilon)</math></b>	<b>3,700M</b> edges
Aggregation/Disaggregation <small>[Chien et al. '04]</small>	$\mathcal{O}(m S  \log 1/\epsilon)$	68M edges
Monte-Carlo <small>[Bahmani et al. '10]</small>	$\mathcal{O}(m + \log m / \epsilon^2)$	68M edges
Power method naive method	$\mathcal{O}(m^2 \log 1/\epsilon)$	11M edges

# Definition of Personalized PageRank

[Brin-Page. Comput. Networks ISDN Syst.'98] [Jeh-Widom. WWW'03]

- Linear equation

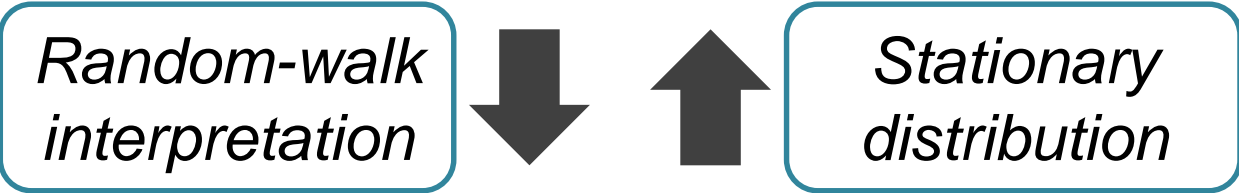
A solution  $x$  of

$$x = \alpha P x + (1 - \alpha) b$$

Preference vector

Transition matrix

Decay factor = 0.85

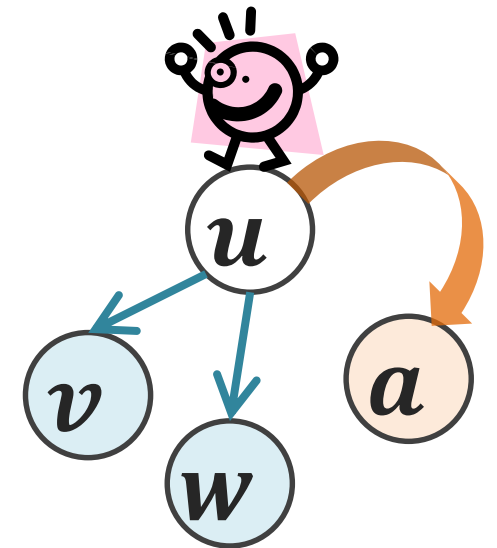


- Random walk modeling web browsing

**Moves** to a random out-neighbor w.p.  $\alpha$

**Jumps** to a random vertex w.p.  $1 - \alpha$

(biased by  $b$ )



# Computing PageRank in Static Graphs

- Solving eq.  $x = \alpha Px + (1 - \alpha)b$

Power method  $x^{(v)} = \alpha Px^{(v-1)} + (1 - \alpha)b$

Gauss-Seidel [Del Corso-Gullí-Romani. Internet Math.'05]

LU/QR factorization [Fujiwara-Nakatsuji-Onizuka-Kitsuregawa. VLDB'12]

Krylov subspace method [Maehara-Akiba-Iwata-Kawarabayashi. VLDB'14]

- Estimating the frequency  $x_v$  of visiting  $v$

Monte-Carlo simulation

# Tracking PageRank in Evolving Graphs

- Aggregation/disaggregation

[Chien-Dwork-Kumar-Simon-Sivakumar. Internet Math.'04]

Apply the power method to a subgraph

↖ Still large 😞

- Monte-Carlo

[Bahmani-Chowdhury. VLDB'10]

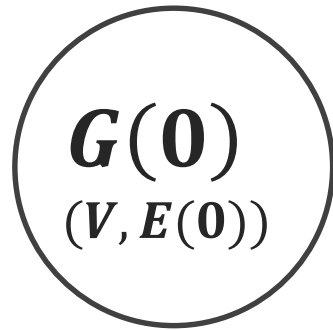
Maintain & update random-walk segments

↖  $\Omega(1/\epsilon^2)$  Too many 😞

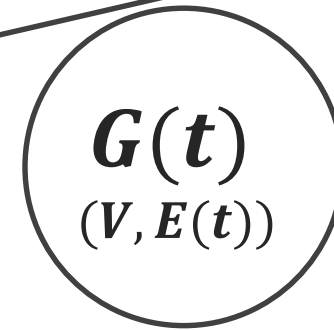
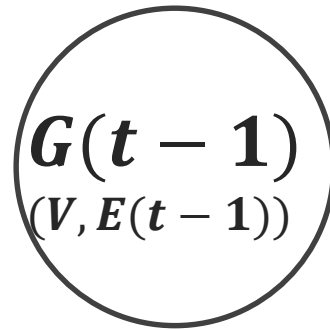
## Proposed Method

# Problem Setting

Given  $G(0), \alpha, b$  at time 0



...



Given at time  $t$ :

Edges **inserted** to / **deleted** from

**Problem at time 0**

Compute **approx.**

PPR  $x(0)$  of  $G(0)$

$$\|x(0) - x^*(0)\|_{\infty} < \epsilon$$

**Problem at time  $t$**

Compute **approx.**

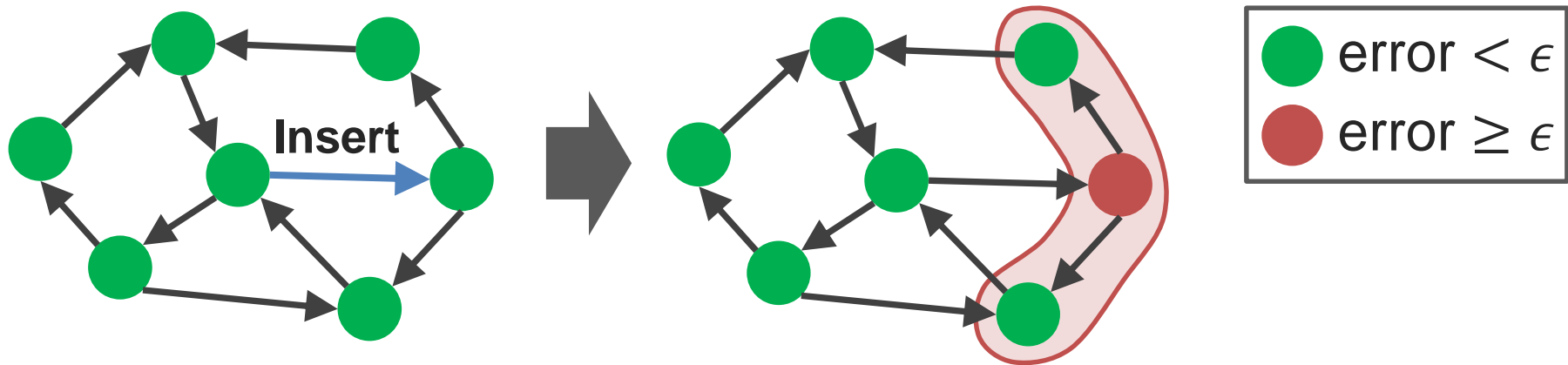
PPR  $x(t)$  of  $G(t)$



## The Idea

Solving eq.  $x(t) = \alpha P(t)x(t) + (1 - \alpha)b$

1.  $x(t - 1)$  is a **GOOD** initial solution for  $x(t)$
2. Improving approximate solution **locally**



- We use the **Gauss-Southwell** method 😊 [Southwell. '40,'46]

*a.k.a. Bookmark coloring algorithm* [Berkhin. Internet Math.'06]

*Local algorithm*

[Spielman-Teng. SIAM J. Comput.'13] [Andersen-Chung-Lang. FOCS'06]

## Proposed Method

# Gauss-Southwell Method [Southwell. '40,'46]

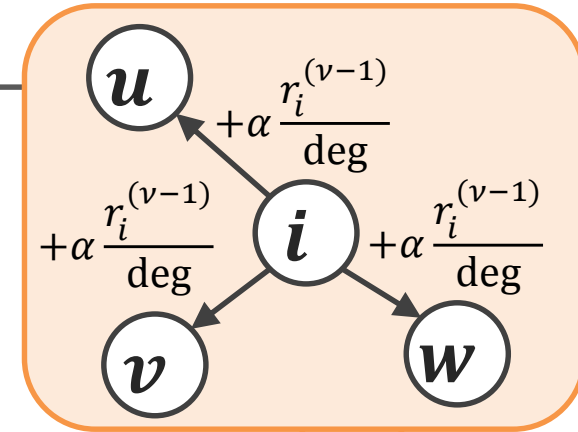
- $\nu^{\text{th}}$  solution  $x^{(\nu)}$
- $\nu^{\text{th}}$  residual  $r^{(\nu)} = (1 - \alpha)b - (I - \alpha P)x^{(\nu)}$   
Goal:  $r^{(\nu)} \rightarrow \mathbf{0}$

$\nu = 1, 2, 3, \dots$

$i \leftarrow$  a vertex with largest  $|r_i^{(\nu-1)}|$

**If**  $|r_i^{(\nu-1)}| < \epsilon$  **terminate**

Update  $x^{(\nu-1)}$  &  $r^{(\nu-1)}$  **locally** so that  $r_i^{(\nu)} = 0$



**# iterations**

Stops within  $\frac{\|r^{(0)}\|}{(1-\alpha)\epsilon}$  iter.  
 $\checkmark \|r^{(\nu)}\| \leq \|r^{(\nu-1)}\| - (1-\alpha)\epsilon$

**Accuracy**

$\|x^* - x^{(\nu)}\|_{\infty} \leq \frac{\epsilon}{1-\alpha}$   
 $\checkmark |r_i^{(\nu)}| < \epsilon$

# Proposed Method Overview

At time  $t$ :

$$x(t)^{(0)} = x(t - 1)$$

$$r(t)^{(0)} = r(t - 1) + \alpha(P(t) - P(t - 1))x(t - 1)$$

Apply the Gauss-Southwell method

# Performance Analysis

At time  $t$ :

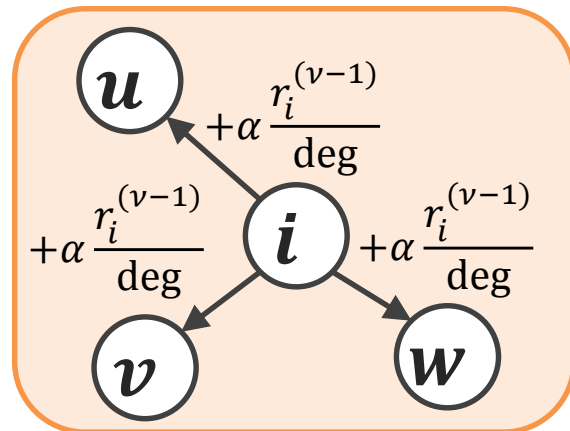
$$x(t)^{(0)} = x(t - 1)$$

$$r(t)^{(0)} = r(t - 1) + \alpha(P(t) - P(t - 1))x(t - 1)$$

Increase of  $\|r\|_1$

Apply the Gauss-Southwell method

Computation time of  $\square = \mathcal{O}(\Delta \times \#iters.)$   
Max. out-deg  $\uparrow$



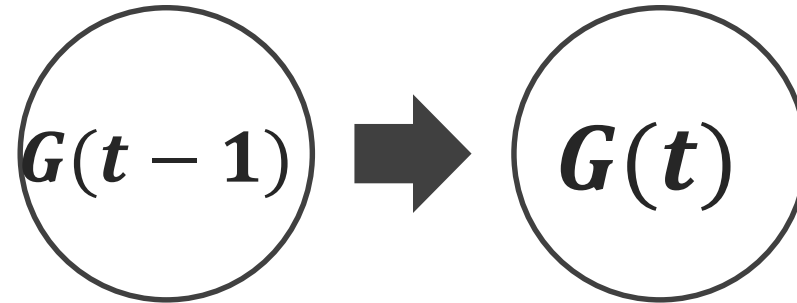
depends on  $\| \dots \|_1$

How small?



# Performance Analysis: Any Change

Consider any change including full construction



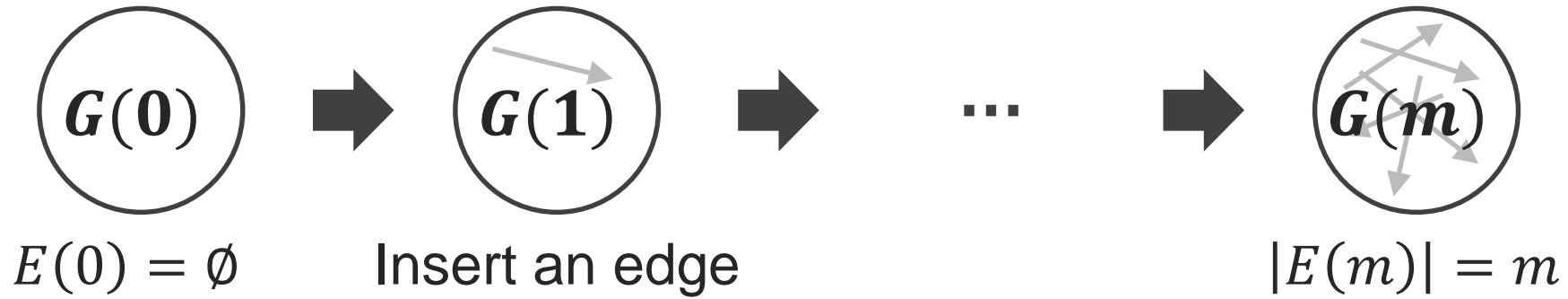
increase of  $\|r\|_1 \leq 2\alpha$

Same as **static computation** 

## Proposed Method

# Performance Analysis: Random Edge Insertion

A single-edge is randomly inserted for each time



**Lemma 3** in [Bahmani-Chowdhury. VLDB'10]

**Monte-Carlo**

[Bahmani-Chowdhury. VLDB'10]

$$\mathbf{E}[\# \text{ updated seg.}] = O(R \log m)$$

$R = \Omega(1/\epsilon^2)$  is total # seg.

**Our method**

$$\mathbf{E}[\text{increase of } \|r\|_1] \leq 2\alpha/t$$

# Performance Analysis: Results

**Our result for random edge insertion** (Prop. 6 in the paper)

If  $m$  edges are randomly and sequentially inserted,  
expected total #iter. is  $\mathcal{O}(\log m / \epsilon)$

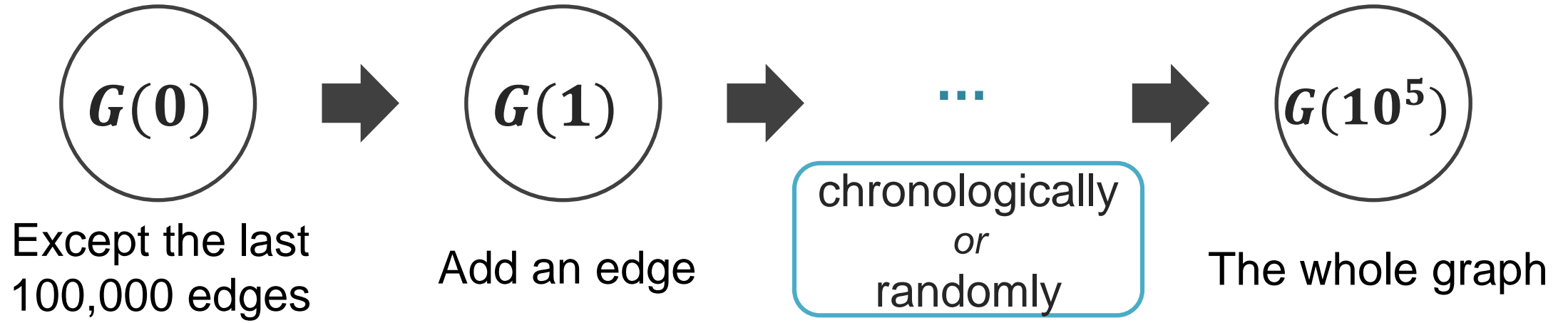
⇒ expected total time is  $\mathcal{O}(m + \Delta \log m / \epsilon)$

**Our result for any change** (Prop. 5 in the paper)

#iter. for any change is amortized  $\mathcal{O}(1/\epsilon)$

⇒ Time is amortized  $\mathcal{O}(\Delta/\epsilon)$

## Setting: Single-edge Insertion



- Parameter settings
  - $\alpha = 0.85$
  - $b$  has 100 non-zero elements
  - $\epsilon = 10^{-9}$



## Experiments

# Efficiency Comparison: Time for an Edge Insertion

	web-Google [SNAP] $ V =1M$ $ E =5M$	Wikipedia [KONECT] $ V =2M$ $ E =40M$	twitter-2010 [LAW] $ V =142M$ $ E =1,500M$	uk-2007-05 [LAW] $ V =105M$ $ E =3,700M$
This work	<b>7</b> $\mu s$	<b>77</b> $\mu s$	<b>29,383</b> $\mu s$	<b>2</b> $\mu s$
Aggregation/Disaggregation [Chien et al. '04]	320 $\mu s$	40,336 $\mu s$	>100,000 $\mu s$	>100,000 $\mu s$
Monte-Carlo [Bahmani et al. '10]	444 $\mu s$	9,196 $\mu s$	>100,000 $\mu s$	>100,000 $\mu s$
Warm start power method	80,994 $\mu s$	>100,000 $\mu s$	>100,000 $\mu s$	>100,000 $\mu s$
From scratch power method	>100,000 $\mu s$	>100,000 $\mu s$	>100,000 $\mu s$	>100,000 $\mu s$

[KONECT] The Koblenz Network Collection <http://konect.uni-koblenz.de/networks/>

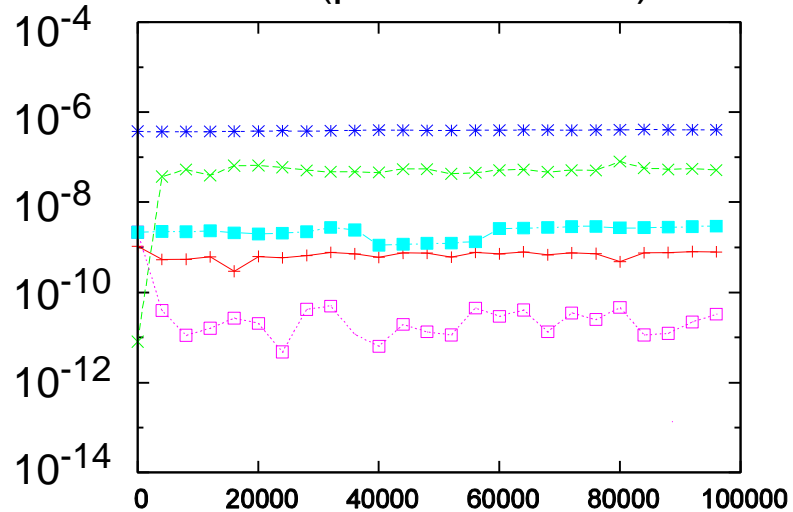
[LAW] Laboratory for Web Algorithmics <http://law.di.unimi.it/datasets.php>

[SNAP] Stanford Large Network Dataset Collection <http://snap.stanford.edu/data/>

# Experiments

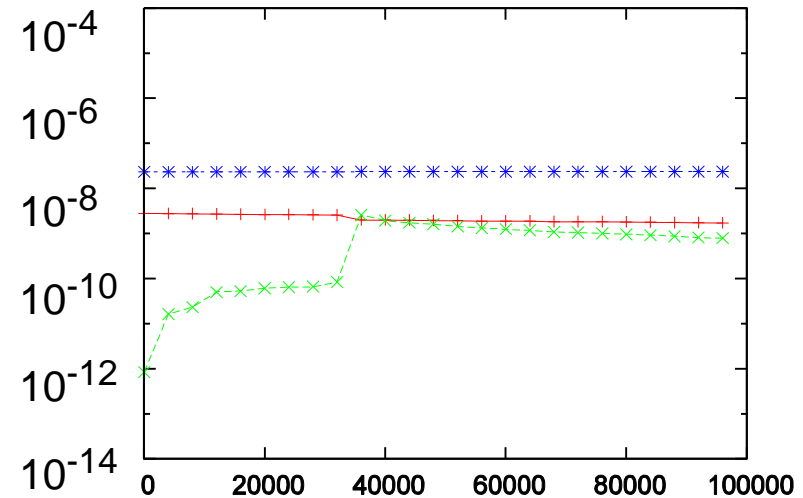
## Accuracy Comparison: Transition of Ave. $L_1$ Error

- +— This work
- - x - - Aggregation/Disaggregation [Chien et al.'04]
- - \* - - Monte-Carlo [Bahmani et al.'10]
- ... □ ... Warm start (power method)
- . ■ - . From scratch (power method)



soc-Epinions1 [SNAP]

$|V|=76K$   $|E|=509K$



Wikipedia [KONECT]

$|V|=2M$   $|E|=40M$

**Comparable** ( $\sim 10^{-9}$ ) to naive methods

Environment: Intel Xeon E5-2690 2.90GHz CPU with 256GB memory

## Experiments

# Evaluation: Time & #Iter. for a Single-edge Insertion

Dataset [Source]	$ V $	$ E $	Max. Out-deg $\Delta$	Ave. Time	Ave. #Iter.
wiki-Talk [SNAP]	2M	5M	100,022	589.6 $\mu$ s	2.3
web-Google [SNAP]	1M	5M	3,444	7.2 $\mu$ s	22.6
as-Skitter [SNAP]	2M	11M	35,387	288.4 $\mu$ s	0.8
Flickr <sup>Time</sup> [KONECT]	2M	33M	26,367	95.3 $\mu$ s	16.2
Wikipedia <sup>Time</sup> [KONECT]	2M	40M	6,975	76.8 $\mu$ s	46.0
soc-LiveJournal1 [SNAP]	5M	68M	20,292	17.9 $\mu$ s	7.6
twitter-2010 [LAW]	42M	1,500M	2,997,469	29,382.8 $\mu$ s	0.7
uk-2007-05 [LAW]	105M	3,700M	15,402	2.3 $\mu$ s	0.0

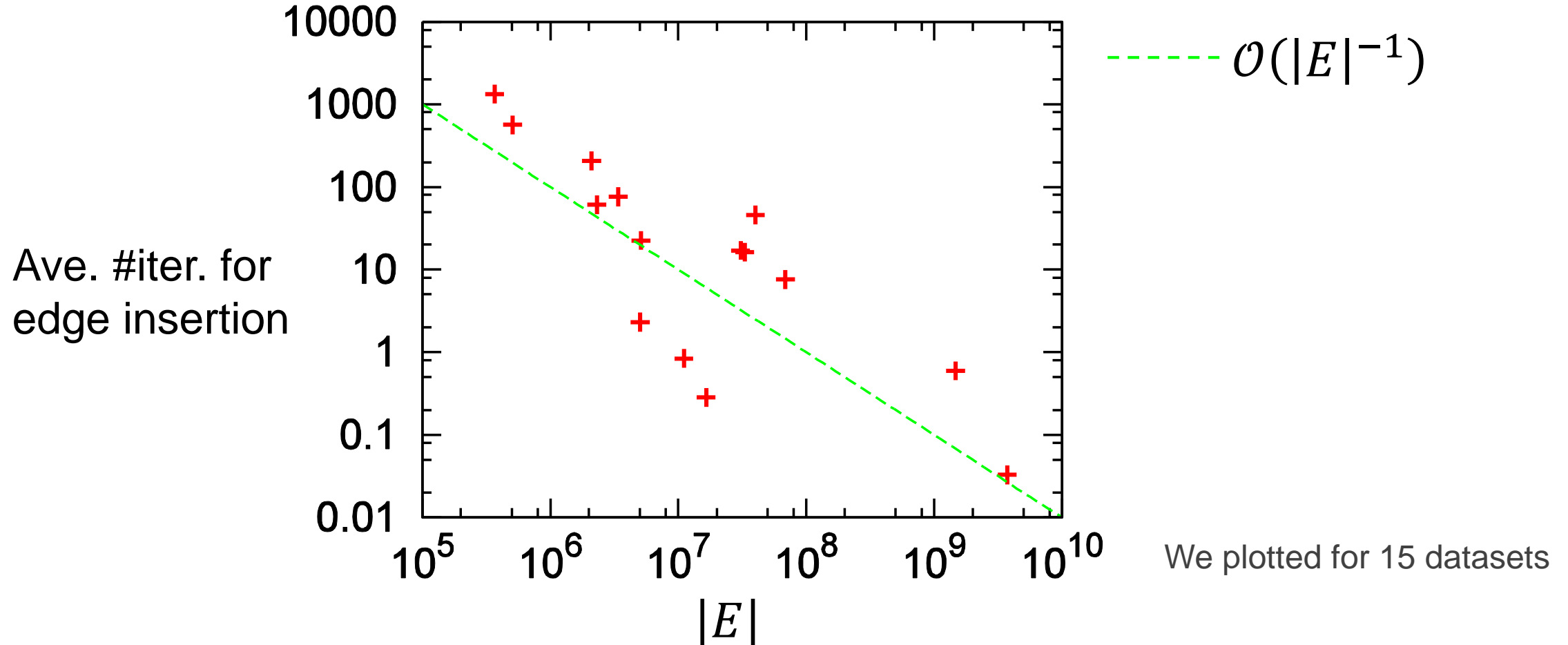
[KONECT] The Koblenz Network Collection <http://konect.uni-koblenz.de/networks/>

[LAW] Laboratory for Web Algorithmics <http://law.di.unimi.it/datasets.php>

[SNAP] Stanford Large Network Dataset Collection <http://snap.stanford.edu/data/>

## Experiments

# Evaluation: Relationship between $|E|$ & #Iter.



Matches our theoretical result

## Experiments

# Evaluation: Time & #Iter. for a Single-edge Insertion

Dataset [Source]	$ V $	$ E $	Max. Out-deg $\Delta$	Ave. Time	Ave. #Iter.
wiki-Talk [SNAP]	2M	5M	100,022	589.6 $\mu$ s	2.3
web-Google [SNAP]	1M	5M	3,444	7.2 $\mu$ s	22.6
as-Skitter [SNAP]	2M	11M	35,387	288.4 $\mu$ s	0.8
Flickr <sup>Time</sup> [KONECT]	2M	33M	26,367	95.3 $\mu$ s	16.2
Wikipedia <sup>Time</sup> [KONECT]	2M	40M	6,975	76.8 $\mu$ s	46.0
soc-LiveJournal1 [SNAP]	5M	68M	20,292	17.9 $\mu$ s	7.6
twitter-2010 [LAW]	42M	1,500M	<b>2,997,469</b>	<b>29,382.8</b> $\mu$ s	0.7
uk-2007-05 [LAW]	105M	3,700M	<b>15,402</b>	<b>2.3</b> $\mu$ s	0.0

[KONECT] The Koblenz Network Collection <http://konect.uni-koblenz.de/networks/>

[LAW] Laboratory for Web Algorithmics <http://law.di.unimi.it/datasets.php>

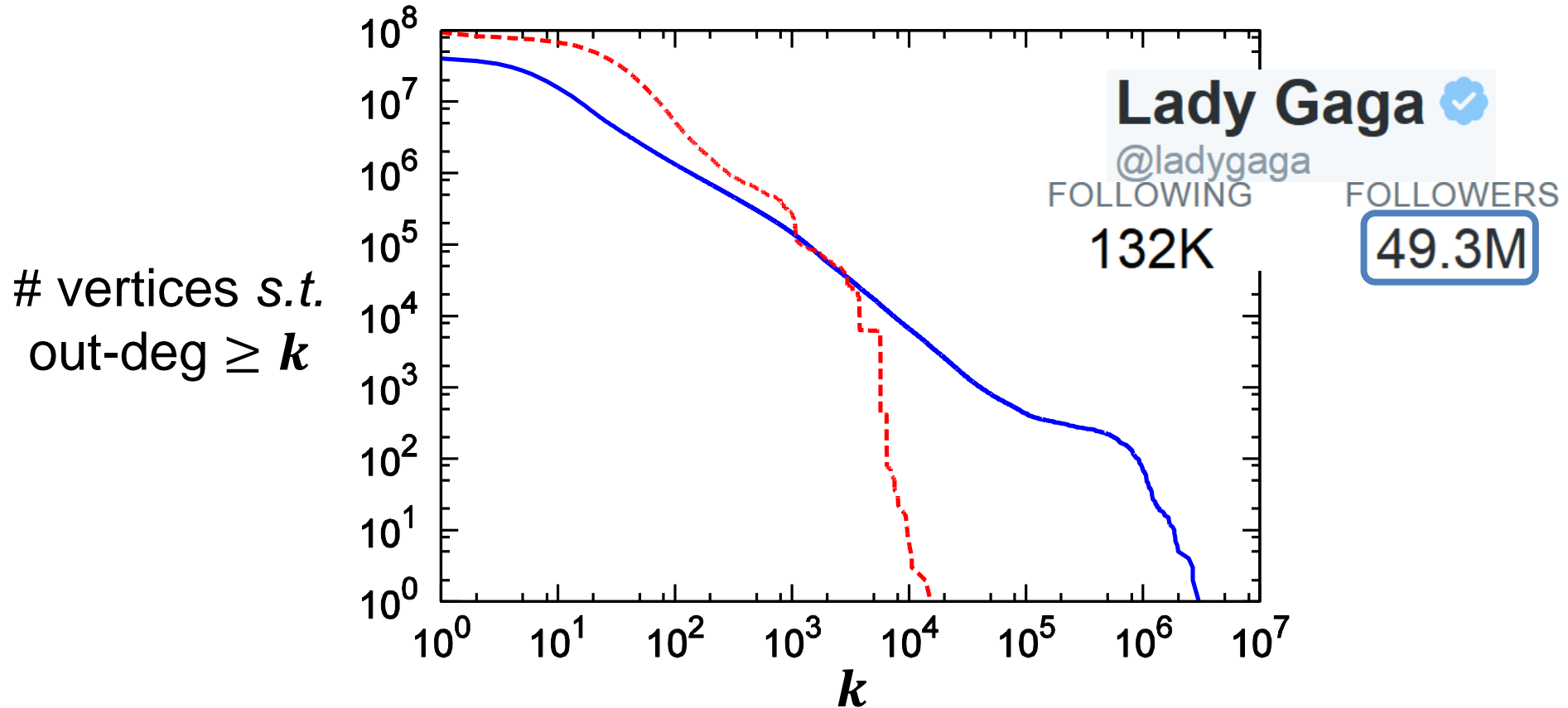
[SNAP] Stanford Large Network Dataset Collection <http://snap.stanford.edu/data/>

**What  
Happens !?**

Environment: Intel Xeon E5-2690 2.90GHz CPU with 256GB memory

# Discussion: What is the Difference?

- twitter-2010  $(u, v)$  says “ $v$  follows  $u$ ”
- - - uk-2007-05  $(u, v)$  says “Page  $u$  links to  $v$ ”



**Celebrities** cause the performance degradation !!

# Summary

Proposed an efficient & accurate method for  
**Personalized PageRank tracking** in evolving networks

## Theoretically

Ave.  $\mathcal{O}(m + \Delta \log m / \epsilon)$   
for  $m$  edge insertions

## Experimentally

Scales to a graph w/ **3.7B edges**

## Future Work

- Further Speed-up based on our observation
- Handle dangling nodes

