

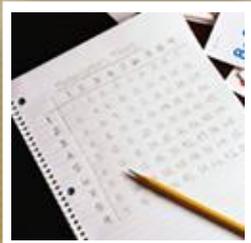
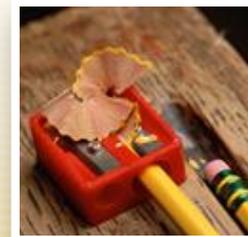


# SimFusion+: Extending SimFusion Towards Efficient Estimation on Large and Dynamic Networks

Weiren Yu<sup>1</sup>, Xuemin Lin<sup>1</sup>, Wenjie Zhang<sup>1</sup>,  
Ying Zhang<sup>1</sup> Jiajin Le<sup>2</sup>,

<sup>1</sup> University of New South Wales & NICTA, Australia

<sup>2</sup> Donghua University, China



# Contents



## 1. Introduction



## 2. Problem Definition



## 3. Optimization Techniques



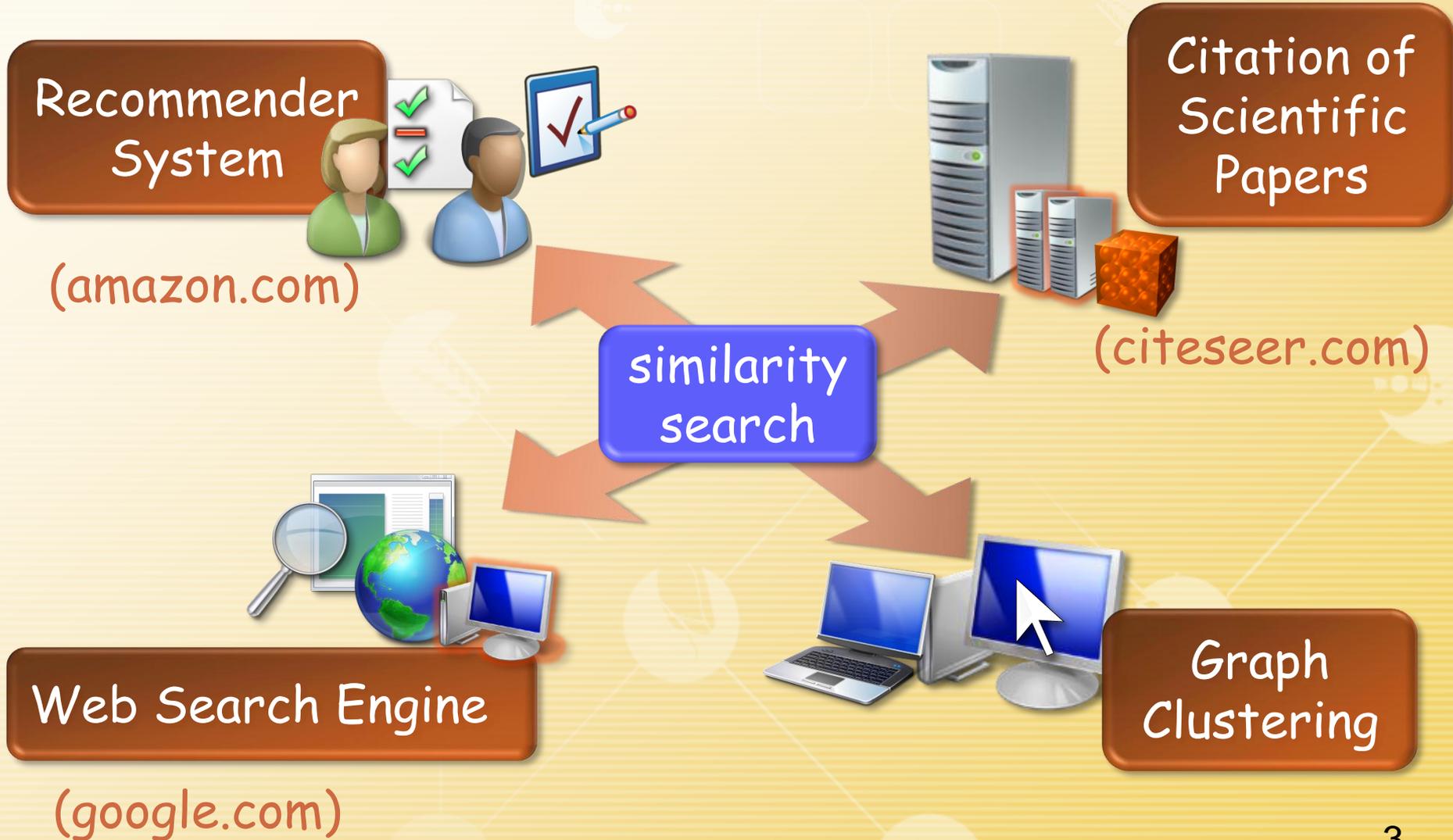
## 4. Experimental Results





# 1. Introduction

- ❖ Many applications require a measure of “similarity” between objects.





# SimFusion: A New Link-based Similarity Measure

## ❖ Structural Similarity Measure

❖ PageRank [Page et. al, 99]

❖ SimRank [Jeh and Widom, KDD 02]

## ❖ SimFusion similarity

❖ A new promising structural measure [Xi et. al, SIGIR 05]

❖ Extension of Co-Citation and Coupling metrics

## ❖ Basic Philosophy

❖ Following **the Reinforcement Assumption**:

*The similarity between objects is reinforced by the similarity of their related objects.*



# SimFusion Overview

## ❖ Features

- ❖ Using a *Unified Relationship Matrix* (URM) to represent relationships among heterogeneous data
- ❖ Defined recursively and is computed iteratively
- ❖ Applicable to any domain with object-to-object relationships

## ❖ Challenges

- ❖ URM may incur trivial solution or divergence issue of SimFusion.
- ❖ Rather costly to compute SimFusion on large graphs
  - ❖ Naïve Iteration: matrix-matrix multiplication
  - ❖ Requiring  $O(Kn^3)$  time,  $O(n^2)$  space [Xi et. al. , SIGIR 05]
- ❖ No incremental algorithms when edges update

# Existing SimFusion: URM and USM

- ❖ *Data Space*:  $\mathcal{D} = \{o_1, o_2, \dots\}$  a finite set of data objects (**vertices**)
- ❖ *Data Relation (edges)* Given an entire space  $\mathcal{D} = \bigcup_{i=1}^N \mathcal{D}_i$ 
  - ❖ *Intra-type Relation*  $\mathcal{R}_{i,i} \subseteq \mathcal{D}_i \times \mathcal{D}_i$  carrying info. within one space
  - ❖ *Inter-type Relation*  $\mathcal{R}_{i,j} \subseteq \mathcal{D}_i \times \mathcal{D}_j$  carrying info. between spaces
- ❖ *Unified Relationship Matrix (URM)*:

$$\mathbf{L}_{\text{URM}} = \begin{pmatrix} \lambda_{1,1} \mathbf{L}_{1,1} & \lambda_{1,2} \mathbf{L}_{1,2} & \cdots & \lambda_{1,N} \mathbf{L}_{1,N} \\ \lambda_{2,1} \mathbf{L}_{2,1} & \lambda_{2,2} \mathbf{L}_{2,2} & \cdots & \lambda_{2,N} \mathbf{L}_{2,N} \\ \vdots & \vdots & \ddots & \vdots \\ \lambda_{N,1} \mathbf{L}_{N,1} & \lambda_{N,2} \mathbf{L}_{N,2} & \cdots & \lambda_{N,N} \mathbf{L}_{N,N} \end{pmatrix} \quad \mathbf{L}_{i,j}(x, y) = \begin{cases} \frac{1}{n_j}, & \text{if } \mathcal{N}_j(x) = \emptyset; \\ \frac{1}{|\mathcal{N}_j(x)|}, & \text{if } (x, y) \in \mathcal{R}_{i,j}; \\ 0, & \text{otherwise.} \end{cases}$$

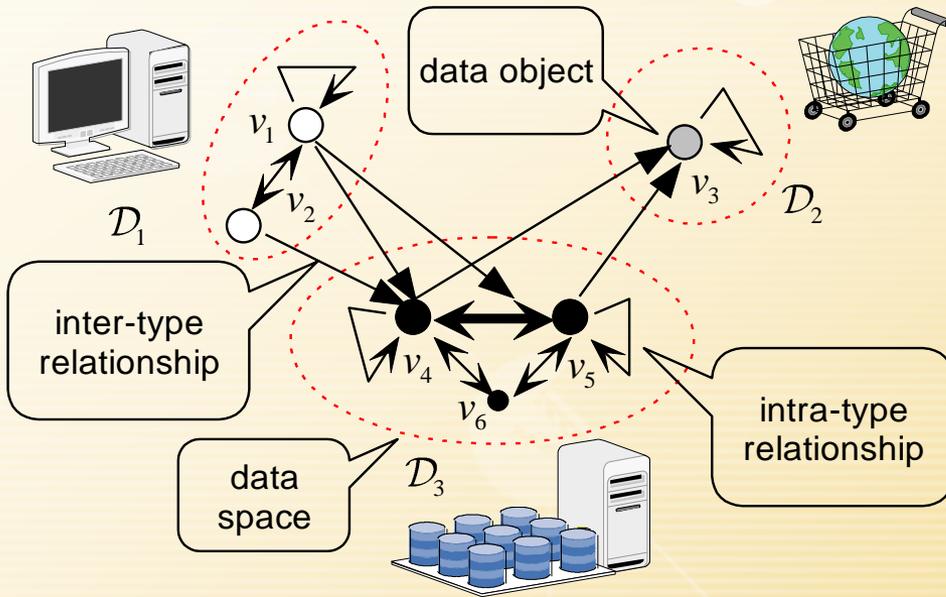
- ❖  $\lambda_{i,j}$  is the *weighting factor* between  $\mathcal{D}_i$  and  $\mathcal{D}_j$
- ❖ *Unified Similarity Matrix (USM)*:

$$\exists \mathbf{S} = \begin{pmatrix} s_{1,1} & \cdots & s_{1,n} \\ \vdots & \ddots & \vdots \\ s_{n,1} & \cdots & s_{n,n} \end{pmatrix} \in \mathbb{R}^{n \times n} \quad \text{s.t.} \quad \mathbf{S} = \mathbf{L} \cdot \mathbf{S} \cdot \mathbf{L}^T.$$



# Example.

## SimFusion Similarity on Heterogeneous Domain



Trivial Solution !!!

$$S = [1]_{n \times n}$$

$$\mathbf{L}_{\text{URM}} = \begin{pmatrix} \frac{1}{8} & \frac{1}{8} & \frac{1}{4} & \frac{1}{4} & \frac{1}{4} & 0 \\ \frac{1}{4} & 0 & \frac{1}{4} & \frac{1}{2} & 0 & 0 \\ \frac{1}{16} & \frac{1}{16} & \frac{1}{4} & \frac{5}{24} & \frac{5}{24} & \frac{5}{24} \\ \frac{1}{10} & \frac{1}{10} & \frac{3}{5} & \frac{1}{15} & \frac{1}{15} & \frac{1}{15} \\ \frac{1}{10} & \frac{1}{10} & \frac{3}{5} & \frac{1}{15} & \frac{1}{15} & \frac{1}{15} \\ \frac{1}{10} & \frac{1}{10} & \frac{3}{5} & \frac{1}{10} & \frac{1}{10} & 0 \end{pmatrix}$$

$$\exists \mathbf{S} \in \mathbb{R}^{n \times n} \quad s.t. \quad \mathbf{S} = \mathbf{L} \cdot \mathbf{S} \cdot \mathbf{L}^T.$$

$$\mathbf{\Lambda} = \begin{matrix} \mathcal{D}_1 & \mathcal{D}_2 & \mathcal{D}_3 \\ \begin{pmatrix} \frac{1}{4} & \frac{1}{4} & \frac{1}{2} \\ \frac{1}{8} & \frac{1}{4} & \frac{5}{8} \\ \frac{1}{5} & \frac{3}{5} & \frac{1}{5} \end{pmatrix} \end{matrix} \quad \mathbf{S}_{\text{USM}} = \begin{pmatrix} 1 & s_{1,2} & \cdots & s_{1,n} \\ s_{2,1} & 1 & \cdots & s_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ s_{n,1} & s_{n,2} & \cdots & 1 \end{pmatrix}$$

High complexity !!!

$O(Kn^3)$  time

$O(n^2)$  space



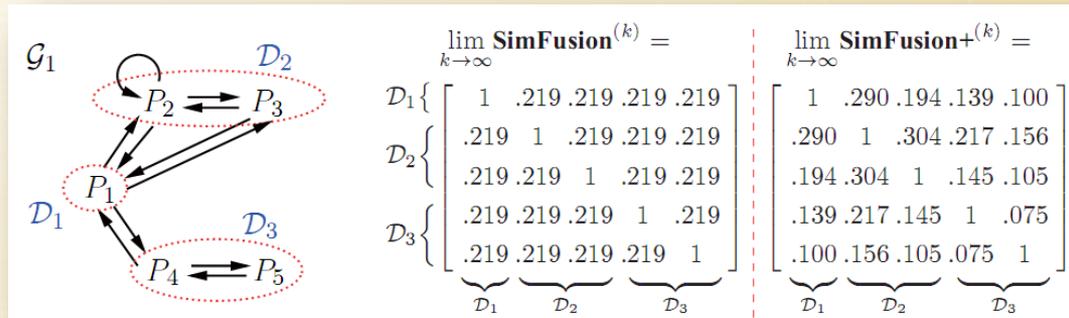
# Contributions

- ❖ Revising the existing SimFusion model, avoiding
  - ❖ non-semantic convergence
  - ❖ divergence issue
- ❖ Optimizing the computation of SimFusion+
  - ❖  $O(Km)$  pre-computation time, plus  $O(1)$  time and  $O(n)$  space
  - ❖ Better accuracy guarantee
- ❖ Incremental computation on edge updates
  - ❖  $O(\delta n)$  time and  $O(n)$  space for handling  $\delta$  edge updates

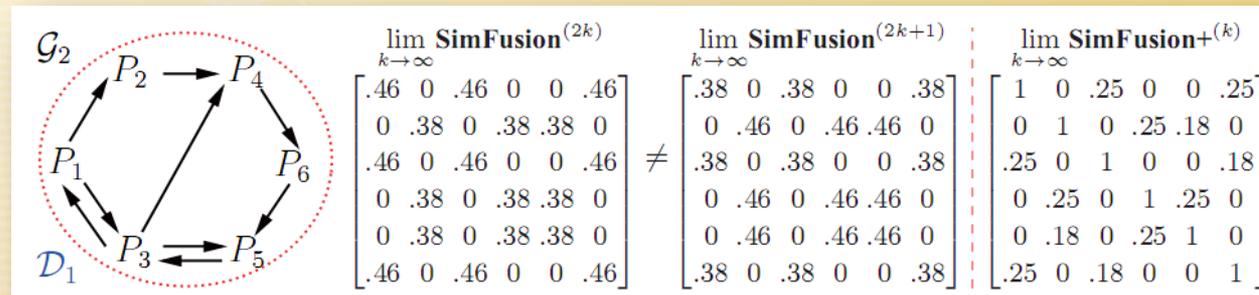
# Revised SimFusion

❖ Motivation: Two issues of the existing SimFusion model

❖ **Trivial** Solution on Heterogeneous Domain



❖ **Divergent** Solution on Homogeneous Domain



**Root cause: row normalization of URM !!!**

# From URM to UAM

❖ Unified Adjacency Matrix (UAM)  $\mathbf{A} = \tilde{\mathbf{A}} + \mathbf{1}/n^2$

$$\tilde{\mathbf{A}} = \begin{pmatrix} \lambda_{1,1}\mathbf{A}_{1,1} & \lambda_{1,2}\mathbf{A}_{1,2} & \cdots & \lambda_{1,N}\mathbf{A}_{1,N} \\ \lambda_{2,1}\mathbf{A}_{2,1} & \lambda_{2,2}\mathbf{A}_{2,2} & \cdots & \lambda_{2,N}\mathbf{A}_{2,N} \\ \vdots & \vdots & \ddots & \vdots \\ \lambda_{N,1}\mathbf{A}_{N,1} & \lambda_{N,2}\mathbf{A}_{N,2} & \cdots & \lambda_{N,N}\mathbf{A}_{N,N} \end{pmatrix}, \quad \mathbf{A}_{i,j}(x,y) = \begin{cases} \frac{1}{n_j}, & \text{if } \mathcal{N}_j(x) = \emptyset; \\ \mathbf{1}, & \text{if } (x,y) \in \mathcal{R}_{i,j}; \\ 0, & \text{otherwise.} \end{cases}$$

❖ Example

$$\mathbf{\Lambda} = \begin{matrix} & \mathcal{D}_1 & \mathcal{D}_2 & \mathcal{D}_3 \\ \mathcal{D}_1 & \begin{bmatrix} \frac{1}{2} \\ \frac{1}{6} \\ \frac{1}{3} \end{bmatrix} & \begin{bmatrix} \frac{1}{6} \\ \frac{7}{12} \\ \frac{1}{4} \end{bmatrix} & \begin{bmatrix} \frac{1}{3} \\ \frac{1}{4} \\ \frac{5}{12} \end{bmatrix} \\ \mathcal{D}_2 & & & \\ \mathcal{D}_3 & & & \end{matrix} \Rightarrow \tilde{\mathbf{A}} = \begin{bmatrix} \frac{1}{2} \begin{bmatrix} 1 \end{bmatrix} & \frac{1}{6} \begin{bmatrix} 1 & 1 \end{bmatrix} & \frac{1}{3} \begin{bmatrix} 1 & 0 \end{bmatrix} \\ \frac{1}{6} \begin{bmatrix} 1 \\ 1 \end{bmatrix} & \frac{7}{12} \begin{bmatrix} 1 & 1 \\ 1 & 0 \end{bmatrix} & \frac{1}{4} \begin{bmatrix} \frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & \frac{1}{2} \end{bmatrix} \\ \frac{1}{3} \begin{bmatrix} 1 \\ 0 \end{bmatrix} & \frac{1}{4} \begin{bmatrix} \frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & \frac{1}{2} \end{bmatrix} & \frac{5}{12} \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \end{bmatrix} = \begin{bmatrix} \frac{1}{2} & \frac{1}{6} & \frac{1}{6} & \frac{1}{3} & 0 \\ \frac{1}{6} & \frac{7}{12} & \frac{7}{12} & \frac{1}{8} & \frac{1}{8} \\ \frac{1}{6} & \frac{7}{12} & 0 & \frac{1}{8} & \frac{1}{8} \\ \frac{1}{3} & \frac{1}{8} & \frac{1}{8} & 0 & \frac{5}{12} \\ 0 & \frac{1}{8} & \frac{1}{8} & \frac{5}{12} & 0 \end{bmatrix}$$

# Revised SimFusion+

## ❖ Basic Intuition

- ❖ replace URM with UAM to postpone “row normalization” in a delayed fashion while preserving the reinforcement assumption of the original SimFusion

## ❖ Revised SimFusion+ Model

$$S = \frac{A \cdot S \cdot A^T}{\|A \cdot S \cdot A^T\|_2}$$

## Original SimFusion

$$S = L \cdot S \cdot L^T$$

squeeze similarity scores in  $S$  into  $[0, 1]$ .



# Optimizing SimFusion+ Computation

## ❖ Conventional Iterative Paradigm

$$\mathbf{S}^{(k+1)} = \frac{\mathbf{A} \cdot \mathbf{S}^{(k)} \cdot \mathbf{A}^T}{\|\mathbf{A} \cdot \mathbf{S}^{(k)} \cdot \mathbf{A}^T\|_2}.$$

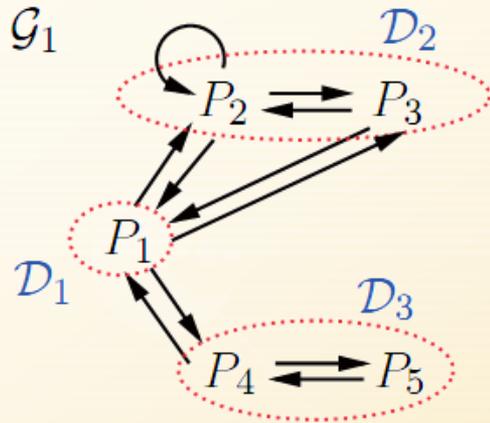
- ❖ Matrix-**matrix** multiplication, requiring  $O(kn^3)$  time and  $O(n^2)$  space
- ❖ Our approach: To convert SimFusion+ computation into finding the dominant eigenvector of the UAM  $\mathbf{A}$ .

$$[\mathbf{S}]_{i,j} = [\sigma_{\max}(\mathbf{A})]_i \times [\sigma_{\max}(\mathbf{A})]_j$$

Pre-compute  $\sigma_{\max}(\mathbf{A})$  only once, and cache it for later reuse

- ❖ Matrix-**vector** multiplication, requiring  $O(km)$  time and  $O(n)$  space

# Example



Assume  $\mathbf{A} = \tilde{\mathbf{A}} + \mathbf{1}/5^2$  with

$$\tilde{\mathbf{A}} = \begin{bmatrix} \frac{1}{2} [1] & \frac{1}{6} [1 \ 1] & \frac{1}{3} [1 \ 0] \\ \frac{1}{6} [1] & \frac{7}{12} [1 \ 1] & \frac{1}{4} \begin{bmatrix} \frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & \frac{1}{2} \end{bmatrix} \\ \frac{1}{3} [1] & \frac{1}{4} \begin{bmatrix} \frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & \frac{1}{2} \end{bmatrix} & \frac{5}{12} [0 \ 1] \end{bmatrix} = \begin{bmatrix} \frac{1}{2} & \frac{1}{6} & \frac{1}{6} & \frac{1}{3} & 0 \\ \frac{1}{6} & \frac{7}{12} & \frac{7}{12} & \frac{1}{8} & \frac{1}{8} \\ \frac{1}{6} & \frac{7}{12} & 0 & \frac{1}{8} & \frac{1}{8} \\ \frac{1}{3} & \frac{1}{8} & \frac{1}{8} & 0 & \frac{5}{12} \\ 0 & \frac{1}{8} & \frac{1}{8} & \frac{5}{12} & 0 \end{bmatrix}$$

## ❖ Conventional Iteration:

$$\mathbf{S}^{(0)} = \mathbf{1} \quad \mathbf{S}^{(k+1)} = \frac{\mathbf{A} \cdot \mathbf{S}^{(k)} \cdot \mathbf{A}^T}{\|\mathbf{A} \cdot \mathbf{S}^{(k)} \cdot \mathbf{A}^T\|_2}$$

$$\mathbf{S} = \begin{bmatrix} .186 & .290 & .194 & .139 & .100 \\ .290 & .453 & .304 & .217 & .156 \\ .194 & .304 & .203 & .145 & .105 \\ .139 & .217 & .145 & .104 & .075 \\ .100 & .156 & .105 & .075 & .054 \end{bmatrix}$$

## ❖ Our approach:

$$\sigma_{\max}(\mathbf{A}) = [.431 \quad .673 \quad .451 \quad .322 \quad .232]^T$$

$$[\mathbf{S}]_{1,2} = [\sigma_{\max}(\mathbf{A})]_1 \times [\sigma_{\max}(\mathbf{A})]_2 = .431 \times .673 = .290.$$

$$[\mathbf{S}]_{1,3} = [\sigma_{\max}(\mathbf{A})]_1 \times [\sigma_{\max}(\mathbf{A})]_3 = .431 \times .451 = .194.$$



# Key Observation

❖ Kronecker product “ $\otimes$ ”:

$$\mathbf{X} \otimes \mathbf{Y} \stackrel{\text{def}}{=} \begin{bmatrix} x_{1,1} \mathbf{Y} & \cdots & x_{1,q} \mathbf{Y} \\ \vdots & \ddots & \vdots \\ x_{p,1} \mathbf{Y} & \cdots & x_{p,q} \mathbf{Y} \end{bmatrix}$$

e.g.  $\mathbf{X} = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}$ ,  $\mathbf{Y} = \begin{bmatrix} 5 & 6 \\ 7 & 8 \end{bmatrix}$ ,  $\mathbf{X} \otimes \mathbf{Y} = \begin{bmatrix} 1 \times \begin{bmatrix} 5 & 6 \\ 7 & 8 \end{bmatrix} & 2 \times \begin{bmatrix} 5 & 6 \\ 7 & 8 \end{bmatrix} \\ 3 \times \begin{bmatrix} 5 & 6 \\ 7 & 8 \end{bmatrix} & 4 \times \begin{bmatrix} 5 & 6 \\ 7 & 8 \end{bmatrix} \end{bmatrix} = \begin{bmatrix} 5 & 6 & 10 & 12 \\ 7 & 8 & 14 & 16 \\ \hline 15 & 18 & 20 & 24 \\ 21 & 24 & 28 & 32 \end{bmatrix}$

❖ Vec operator:  $\text{vec}(\mathbf{X}) \stackrel{\text{def}}{=} [x_{1,1}, \cdots, x_{p,1}, \cdots, x_{1,q}, \cdots, x_{p,q}]^T$

e.g.  $\text{vec}(\mathbf{X}) = [1 \ 3 \ 2 \ 4]^T$

❖ Two important Properties:

$$\text{vec}(\mathbf{BCD}^T) = (\mathbf{D} \otimes \mathbf{B}) \cdot \text{vec}(\mathbf{C})$$

$$\sigma_{\max}(\mathbf{A} \otimes \mathbf{A}) = \sigma_{\max}(\mathbf{A}) \otimes \sigma_{\max}(\mathbf{A})$$



# Key Observation

❖ Two important Properties:

$$\text{P1. } \text{vec}(\mathbf{BCD}^T) = (\mathbf{D} \otimes \mathbf{B}) \cdot \text{vec}(\mathbf{C})$$

$$\text{P2. } \sigma_{\max}(\mathbf{A} \otimes \mathbf{A}) = \sigma_{\max}(\mathbf{A}) \otimes \sigma_{\max}(\mathbf{A})$$

❖ Our main idea:

$$\begin{aligned} \mathbf{S}^{(k+1)} &= \frac{\mathbf{A} \cdot \mathbf{S}^{(k)} \cdot \mathbf{A}^T}{\|\mathbf{A} \cdot \mathbf{S}^{(k)} \cdot \mathbf{A}^T\|_2} \xrightarrow{(1)} \text{vec}(\mathbf{S}^{(k+1)}) = \frac{(\mathbf{A} \otimes \mathbf{A}) \cdot \text{vec}(\mathbf{S}^{(k)})}{\|(\mathbf{A} \otimes \mathbf{A}) \cdot \text{vec}(\mathbf{S}^{(k)})\|_2} \\ &\quad \downarrow \text{Power Iteration} \\ \text{vec}(\mathbf{S}) &= \sigma_{\max}(\mathbf{A}) \otimes \sigma_{\max}(\mathbf{A}) \xleftarrow{(2)} \lim_{k \rightarrow \infty} \text{vec}(\mathbf{S}^{(k)}) = \sigma_{\max}(\mathbf{A} \otimes \mathbf{A}) \end{aligned}$$

# Accuracy Guarantee

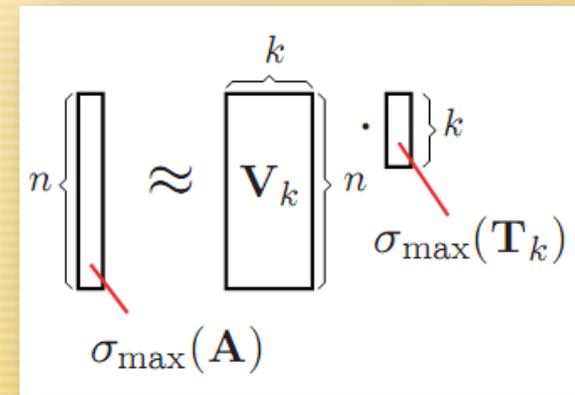
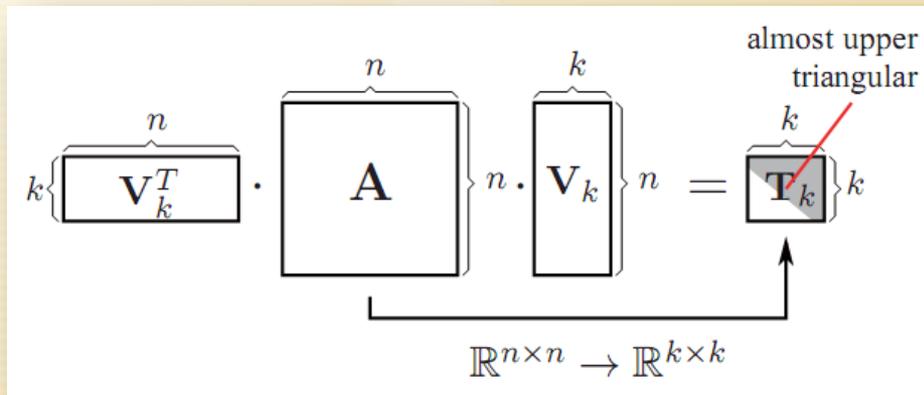
- ❖ Conventional Iterations: **No accuracy guarantee !!!**

$$S^{(k+1)} = \frac{A \cdot S^{(k)} \cdot A^T}{\|A \cdot S^{(k)} \cdot A^T\|_2}$$

$$S = \frac{A \cdot S \cdot A^T}{\|A \cdot S \cdot A^T\|_2}$$

Question:  $\|S^{(k+1)} - S\| \leq ?$

- ❖ Our Method: Utilize Arnoldi decomposition to build an order- $k$  orthogonal subspace for the UAM  $A$ .



Due to  $T_k$  small size and almost "upper-triangularity", Computing  $\sigma_{\max}(T_k)$  is less costly than  $\sigma_{\max}(A)$ .



# Accuracy Guarantee

## ❖ Arnoldi Decomposition:

$$\mathbf{V}_k^T \mathbf{A} \mathbf{V}_k = \mathbf{T}_k,$$

$$\mathbf{V}_k = [\mathbf{v}_1 \mid \mathbf{v}_2 \mid \cdots \mid \mathbf{v}_k]$$

$$\mathbf{A} \mathbf{V}_k - \mathbf{V}_k \mathbf{T}_k = \delta_k \mathbf{v}_{k+1} \mathbf{e}_k^T,$$

$$\mathbf{T}_{k+1} = \begin{bmatrix} \mathbf{T}_k & \star \\ \star & \star \end{bmatrix}$$

## ❖ k-th iterative similarity

$$[\hat{\mathbf{S}}_k]_{i,j} = [\mathbf{V}_k \cdot \sigma_{\max}(\mathbf{T}_k)]_i \times [\mathbf{V}_k \cdot \sigma_{\max}(\mathbf{T}_k)]_j$$

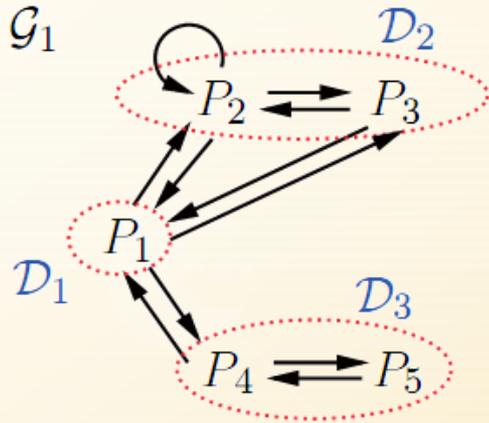
## ❖ Estimate Error:

$$\|\hat{\mathbf{S}}_k - \mathbf{S}\|_2 \leq \epsilon_k$$

$$\epsilon_k = 2 \times |\delta_k \times [\sigma_{\max}(\mathbf{T}_k)]_k|$$



# Example



Assume  $\mathbf{A} = \tilde{\mathbf{A}} + 1/5^2$  with

$$\tilde{\mathbf{A}} = \begin{bmatrix} \frac{1}{2} & \frac{1}{6} & \frac{1}{6} & \frac{1}{3} & 0 \\ \frac{1}{6} & \frac{7}{12} & \frac{7}{12} & \frac{1}{8} & \frac{1}{8} \\ \frac{1}{6} & \frac{7}{12} & 0 & \frac{1}{8} & \frac{1}{8} \\ \frac{1}{3} & \frac{1}{8} & \frac{1}{8} & 0 & \frac{5}{12} \\ 0 & \frac{1}{8} & \frac{1}{8} & \frac{5}{12} & 0 \end{bmatrix}$$

Given  $\epsilon = 0.05$

$$\epsilon_k = 2 \times |\delta_k \times [\sigma_{\max}(\mathbf{T}_k)]_k|$$

❖ Arnoldi Decomposition:

$k$	$\mathbf{T}_k$	$\mathbf{v}_{k+1}$	$\delta_k$	$\sigma_{\max}(\mathbf{T}_k)$	$\epsilon_k$
0	—	$[\.447 \ .447 \ .447 \ .447 \ .447]^T$	—	—	—
1	$[1.08]$	$[\.125 \ .750 \ -.125 \ -.125 \ -.625]^T$	.298	$[1]$	.596
2	$\begin{bmatrix} 1.08 & .298 \\ .298 & .190 \end{bmatrix}$	$[-.089 \ .044 \ .710 \ -.697 \ .032]^T$	.359	$\begin{bmatrix} .957 \\ .290 \end{bmatrix}$	.208
3	$\begin{bmatrix} 1.08 & .298 & 0 \\ .298 & .190 & .359 \\ 0 & .359 & -.083 \end{bmatrix}$	$[-.881 \ .329 \ .137 \ .280 \ .135 \ .231]^T$	.231	$\begin{bmatrix} .945 \\ .316 \\ .090 \end{bmatrix}$	.041

(2)

$$\hat{\mathbf{S}}_3 = \begin{bmatrix} .206 & .301 & .203 & .146 & .103 \\ .301 & .440 & .296 & .213 & .151 \\ .203 & .296 & .199 & .143 & .102 \\ .146 & .213 & .143 & .102 & .073 \\ .103 & .151 & .102 & .073 & .051 \end{bmatrix}$$

(3)

$$\mathbf{A}\mathbf{V}_k - \mathbf{V}_k\mathbf{T}_k = \delta_k \mathbf{v}_{k+1} \mathbf{e}_k^T,$$

$$[\hat{\mathbf{S}}_k]_{i,j} = [\mathbf{V}_k \cdot \sigma_{\max}(\mathbf{T}_k)]_i \times [\mathbf{V}_k \cdot \sigma_{\max}(\mathbf{T}_k)]_j$$

# Edge Update on Dynamic Graphs

## ❖ Incremental UAM

Given old  $G=(D,R)$  and a new  $G'=(D,R')$ , the incremental UAM is a list of edge updates, i.e.,  $\bar{A} = A' - A$

## ❖ Main idea

To reuse  $\bar{A}$  and the eigen-pair  $(\alpha_p, \xi_p)$  of the old  $A$  to compute  $S'$   
 $\bar{A}$  is a sparse matrix when the number  $\delta$  of edge updates is small.

## ❖ Incrementally computing SimFusion+

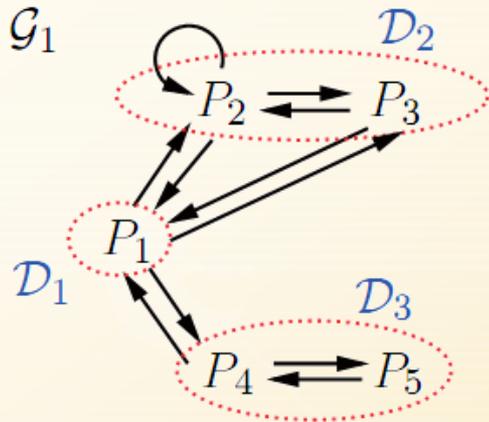
$$[S']_{i,j} = [\xi']_i \cdot [\xi']_j \text{ with } [\xi']_i = [\xi_1]_i + \sum_{p=2}^n c_p \times [\xi_p]_i$$

$$c_p = \frac{\xi_p^T \cdot \eta}{\alpha_p - \alpha_1} \text{ and } \eta = \bar{A} \cdot \xi_1$$

$O(\delta n)$  time

$O(n)$  space

# Example



Suppose edges (P1,P2) and (P2,P1) are removed.

$$\bar{\mathbf{A}} = \begin{bmatrix} 0 & -\frac{1}{6} & 0 & 0 & 0 \\ -\frac{1}{6} & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

$p$	$\alpha_p$	$\xi_p$	$c_p$
1	1.184	$[\.431 \ .673 \ .451 \ .322 \ .232]^T$	—
2	.503	$[\.708 \ -.522 \ -.242 \ .388 \ .132]^T$	.062
3	-.480	$[-.256 \ -.020 \ .095 \ .716 \ -.641]^T$	-.018
4	-.366	$[-.021 \ -.507 \ .853 \ -.119 \ .017]^T$	-.025
5	.242	$[\.497 \ .127 \ .037 \ -.467 \ -.719]^T$	.069

$$\eta = \bar{\mathbf{A}} \cdot \xi_1 = [-.112 \quad -.072 \quad 0 \quad 0 \quad 0]^T.$$

$$c_2 = \xi_2^T \cdot \eta / (\alpha_2 - \alpha_1) = -.0419 / (.503 - 1.184) = .062,$$

$$c_3 = \xi_3^T \cdot \eta / (\alpha_3 - \alpha_1) = .030 / (-.480 - 1.184) = -.018.$$

$$\xi' = \xi_1 + \sum_{p=2}^5 c_p \times \xi_p = [.327 \ .703 \ .485 \ .326 \ .266]^T$$

$$\mathbf{S}' = \begin{bmatrix} .107 & .230 & .159 & .107 & .087 \\ .230 & .494 & .341 & .230 & .187 \\ .159 & .341 & .235 & .158 & .129 \\ .107 & .230 & .158 & .107 & .087 \\ .087 & .187 & .129 & .087 & .071 \end{bmatrix}$$



# Experimental Setting

## ❖ Datasets

- ❖ Synthetic data (RAND 0.5M-3.5M)
- ❖ Real data (DBLP, WEBKB)

DBLP

	$\mathcal{G}_1: 01-02$	$\mathcal{G}_2: 01-04$	$\mathcal{G}_3: 01-06$	$\mathcal{G}_4: 01-08$	$\mathcal{G}_5: 01-10$
$ \mathcal{D} $	1,838	3,723	5,772	9,567	12,276
$ \mathcal{R} $	7,103	14,419	29,054	45,310	64,208

WEBKB

	$U_1: CO$	$U_2: TE$	$U_3: WA$	$U_4: WI$
$ \mathcal{D} $	867	827	1,263	1,205
$ \mathcal{R} $	1,496	1,428	2,969	1,805

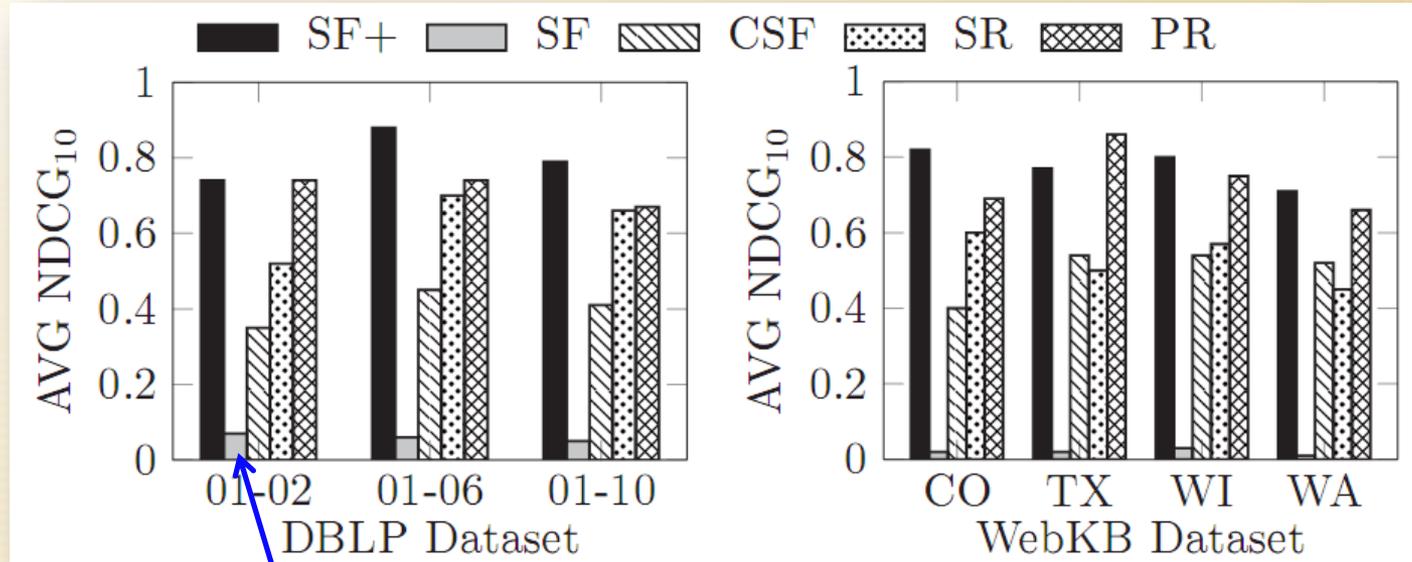
## ❖ Compared Algorithms

- ❖ **SimFusion+** and **IncSimFusion+** ;
- ❖ **SF**, a SimFusion algorithm via matrix iteration [Xi et. al, SIGIR 05];
- ❖ **CSF**, a variant SF, using PageRank distribution [Cai et. al, SIGIR 10];
- ❖ **SR**, a SimRank algorithm via partial sums [Lizorkin et. al, VLDBJ 10];
- ❖ **PR**, a P-Rank encoding both in- and out-links [Zhao et. al, CIKM 09];

# Experiment (1): Accuracy

$$\text{NDCG}_p = \frac{1}{\text{IDCG}_p} \sum_{i=1}^p (2^{\text{rank}_i} - 1) / (\log_2(1 + i))$$

On DBLP and WEBKB



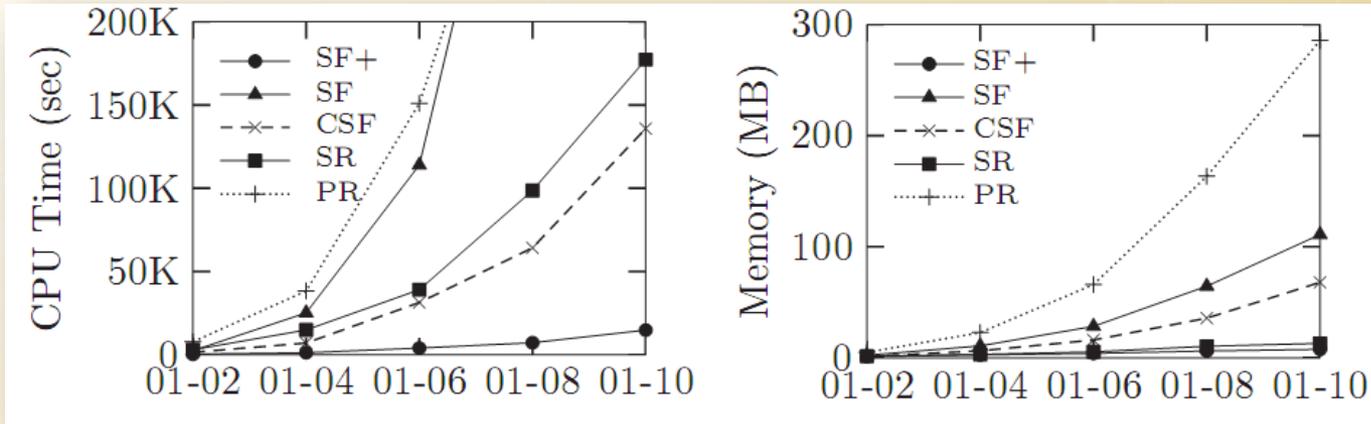
SF+ accuracy is consistently stable on different datasets.

SF seems hardly to get sensible similarities as all its similarities asymptotically approach the same value as K grows.



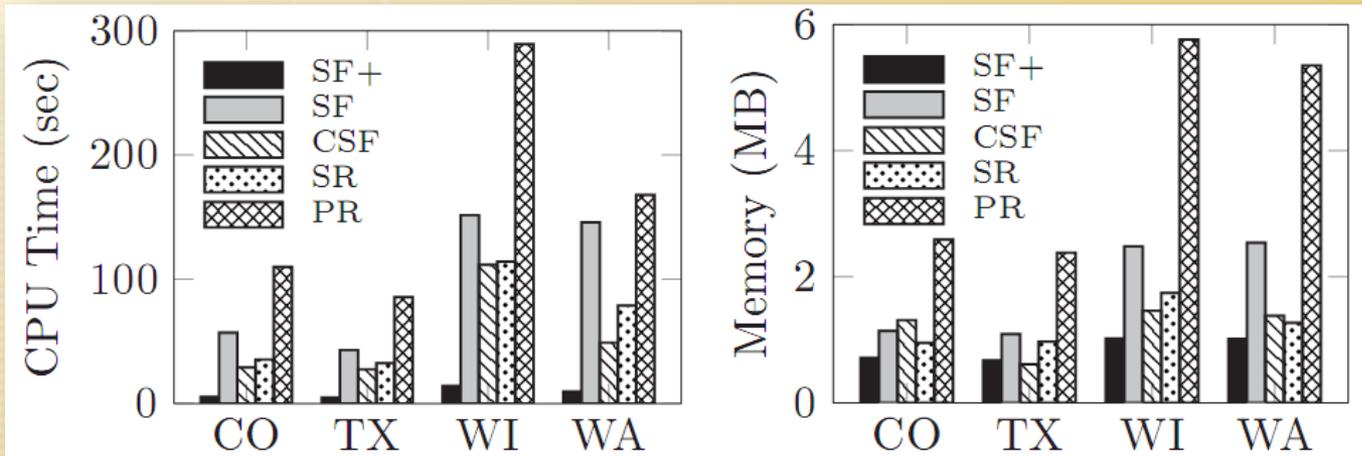
# Experiment (2): CPU Time and Space

On DBLP



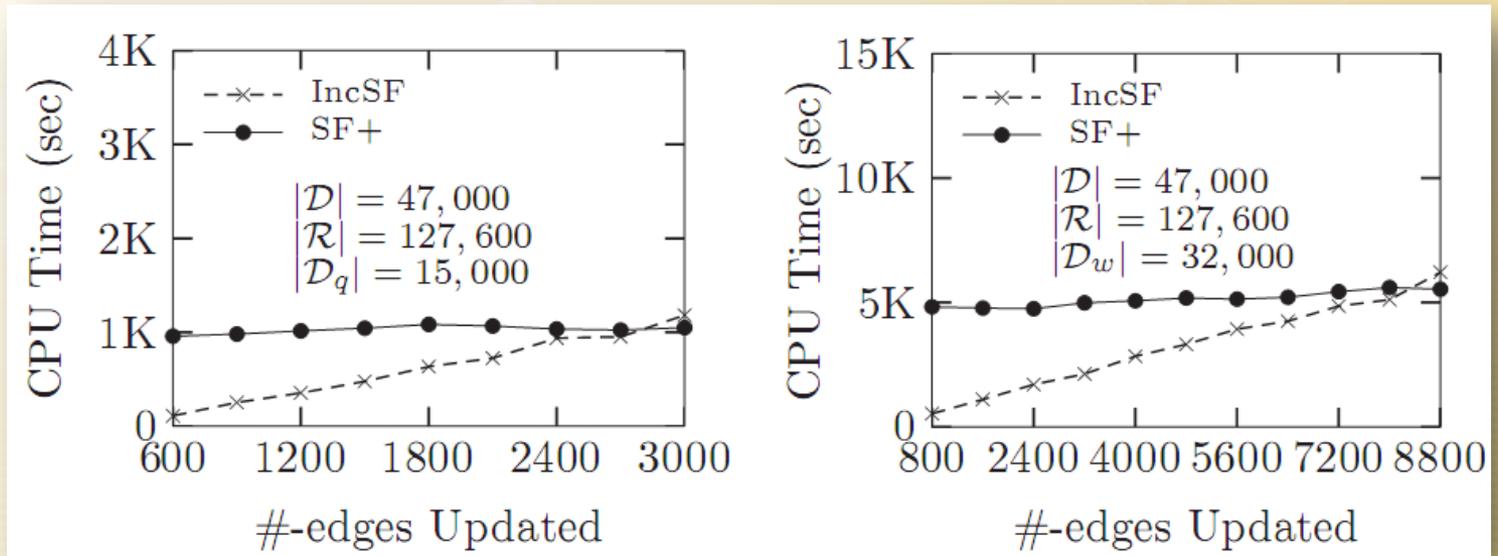
SF+ outperforms the other approaches, due to the use of  $\sigma_{\max}(T_k)$

On WEBKB



# Experiment (3): Edge Updates

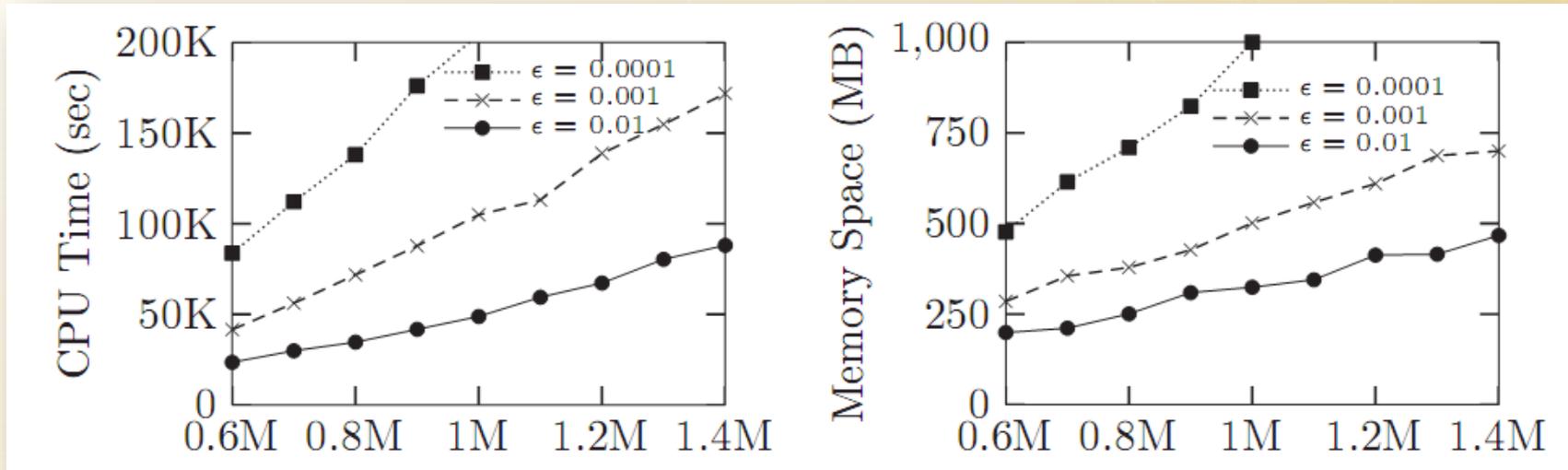
Varying  $\delta$



IncSF+ outperformed SF+ when  $\delta$  is small.

For larger  $\delta$ , IncSF+ is not that good because the small value of  $\delta$  preserves the sparseness of the incremental UAM.

## Experiment (4) : Effects of $\epsilon$



The small choice of  $\epsilon$  imposes more iterations on computing  $T_k$  and  $v_k$ , and hence increases the estimation costs.



## Conclusions

- ❖ A revision of SimFusion+, for preventing the trivial solution and the divergence issue of the original model.
- ❖ Efficient techniques to improve the time and space of SimFusion+ with accuracy guarantees.
- ❖ An incremental algorithm to compute SimFusion+ on dynamic graphs when edges are updated.

## Future Work

- ❖ Devise vertex-updating methods for incrementally computing SimFusion+.
- ❖ Extend to parallelize SimFusion+ computing on GPU.



**Thank You !**