

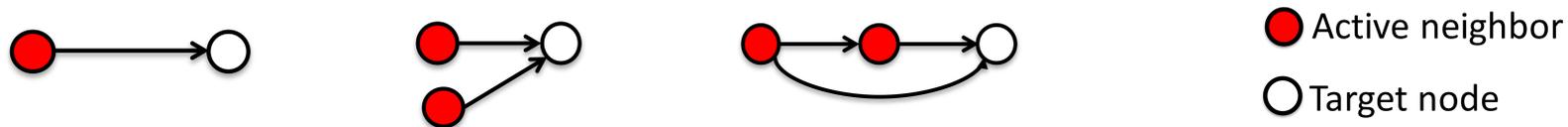


# StructInf: Mining Structural Influence from Social Streams

Jing Zhang\*, Jie Tang<sup>+</sup>, Yuanyi Zhong<sup>+</sup>, Yuchen Mo<sup>+</sup>, Juanzi Li<sup>+</sup>,  
Guojie Song<sup>#</sup>, Wendy Hall<sup>°</sup>, and Jimeng Sun<sup>△</sup>

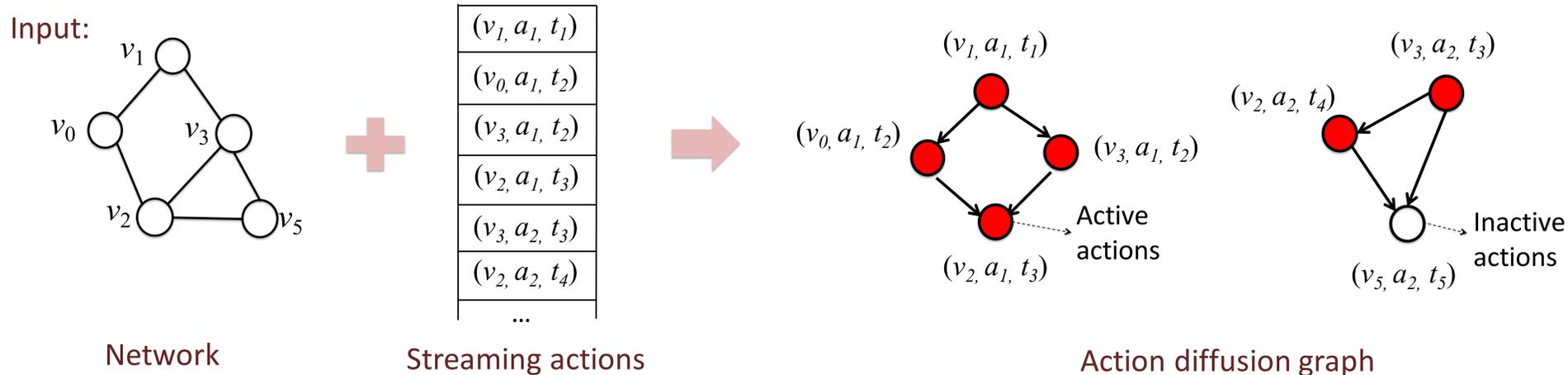
\*Renmin University of China, <sup>+</sup>Tsinghua University,

<sup>#</sup>Peking University, <sup>°</sup>University of Southampton, <sup>△</sup> Georgia Institute of Technology



Question: In which structures, the target nodes are most likely to be activated?

## Problem Formulation



Output:

Influence Probability		$l_t \in L$	$l_t \notin L$
$IP(C_k) = \frac{x_k}{x_k + y_k}$	$\frac{C_k}{C_k}$	$x_k$	$y_k$
		$z_k$	$w_k$

**Structural influence:**  
Influence Probabilities of a structure  $C_k$  that can be found in the action diffusion graph.

## Structural Influence Measurement

### Basic method

- Maintain a queue and a map to record the diffusion edges within recent time interval.
- To calculate  $x_k$ , active actions are newly arrived actions.
- To calculate  $y_k$ , inactive actions are actions that are outdated.
- Enumerate structures by extending neighboring actions of active or inactive actions.
- To avoid duplicate enumeration, assign each action an incremental (unique) label when it arrives, and make the labels of the selected actions smaller than those in the candidate actions.

### Sampling based methods

- Sampling1 – Randomly sample nodes when enumerating influence patterns.
- Sampling2 – Randomly reserve edges when building diffusion graph.
- Sampling3 – Combine Sampling1 and Sampling2.

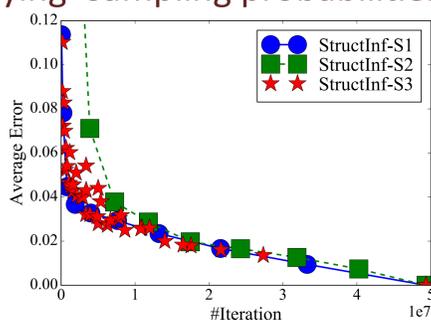
They are unbiased sampling methods

## Results

$k$	$C_k$	$IP_k$	$\tilde{IP}_k$	$U_{IP_k}$	$k$	$C_k$	$IP_k$	$\tilde{IP}_k$	$U_{IP_k}$
1		0.066	0.066	0.020	11		0.038	0.038	0.720
2		0.074	0.074	0.085	12		0.186	0.186	0.088
3		0.111	0.110	0.425	13		0.399	0.392	1.785
4		0.307	0.304	0.928	14		0.063	0.062	0.616
5		0.069	0.069	0.530	15		0.619	<b>0.615</b>	0.548
6		0.091	0.090	0.358	16		0.444	<b>0.439</b>	1.378
7		0.067	0.067	0.236	17		0.070	0.070	0.074
8		0.106	0.099	5.852	18		0.420	<b>0.416</b>	0.890
9		0.381	0.388	1.666	19		0.662	<b>0.645</b>	2.696
10		0.165	0.162	1.128	20		0.485	<b>0.479</b>	1.239

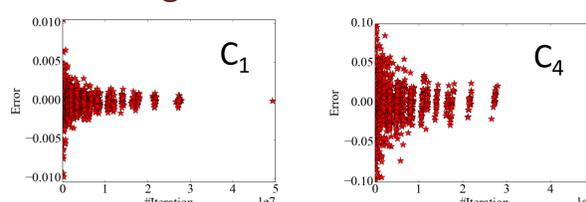
$\tilde{IP}_k$ : Approximate influence probability  
 $U_{IP_k}$ : Relative error of approximate values

Trade-off between error and time by varying sampling probabilities:

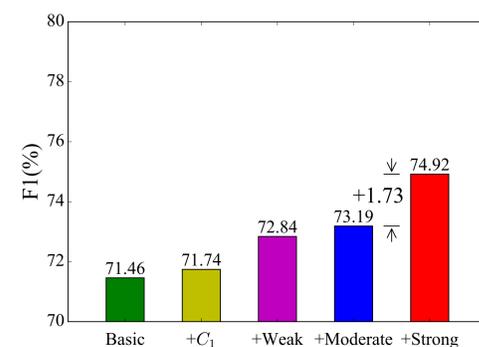


Sampling3 is most insensitive to parameters

Convergence of relative error:



Retweet prediction:



Basic: #friends, gender, status, etc.

$C_1$ : the number of active neighbors

Weak:  $\tilde{IP}_k < 0.1$

Moderate:  $0.1 \leq \tilde{IP}_k < 0.3$

Strong:  $\tilde{IP}_k > 0.3$