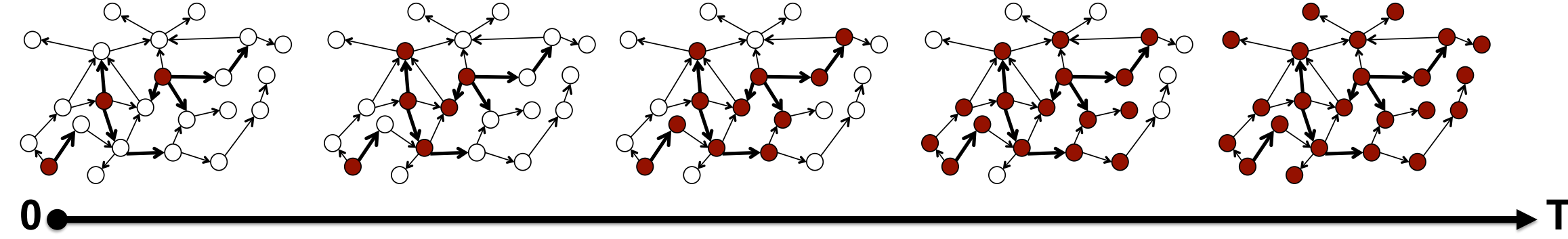


## MOTIVATION

- Question : How can we optimize the selection of the earlier nodes to trigger, *within a time window  $T$* , the largest expected number of follow-ups ?

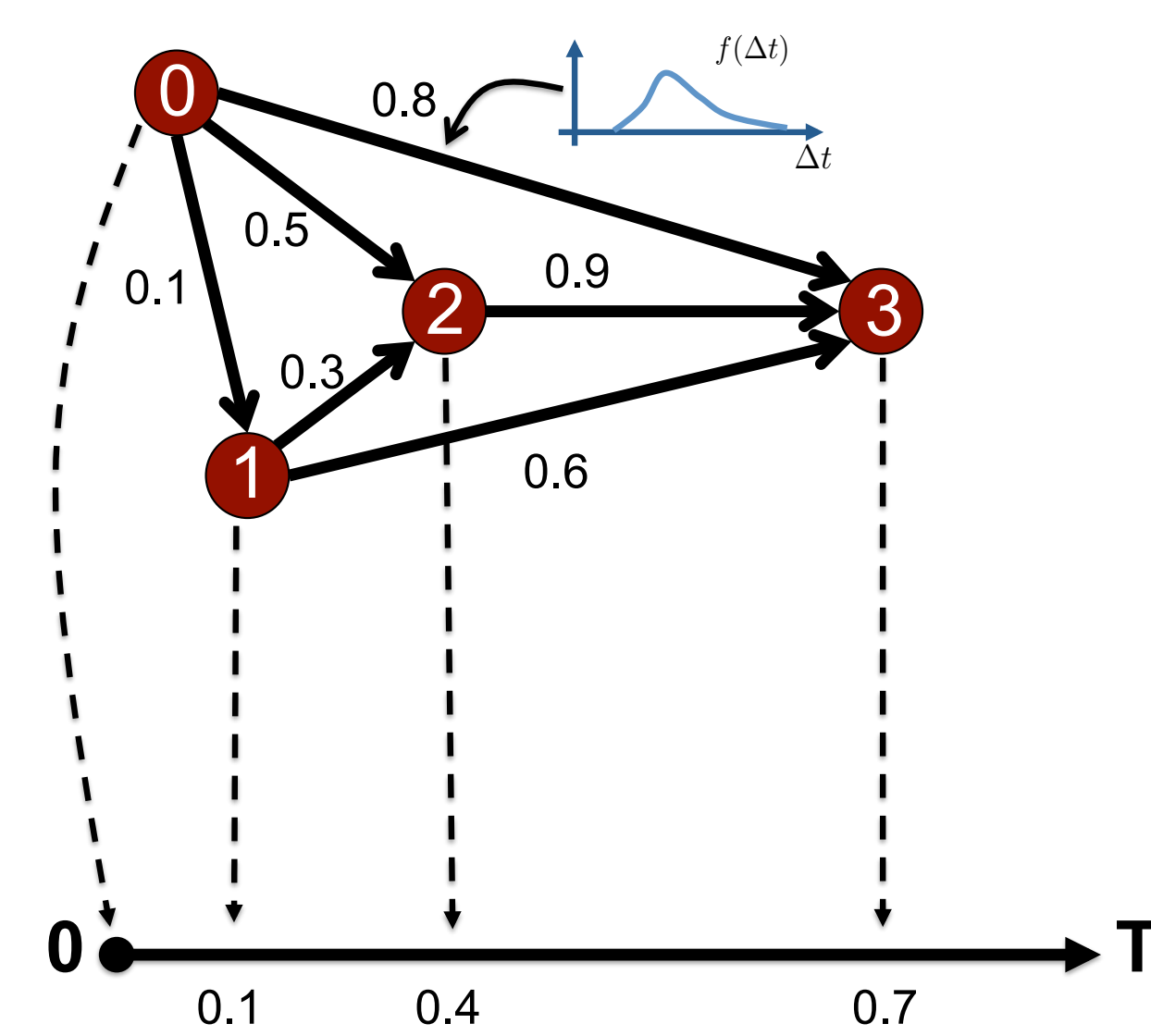
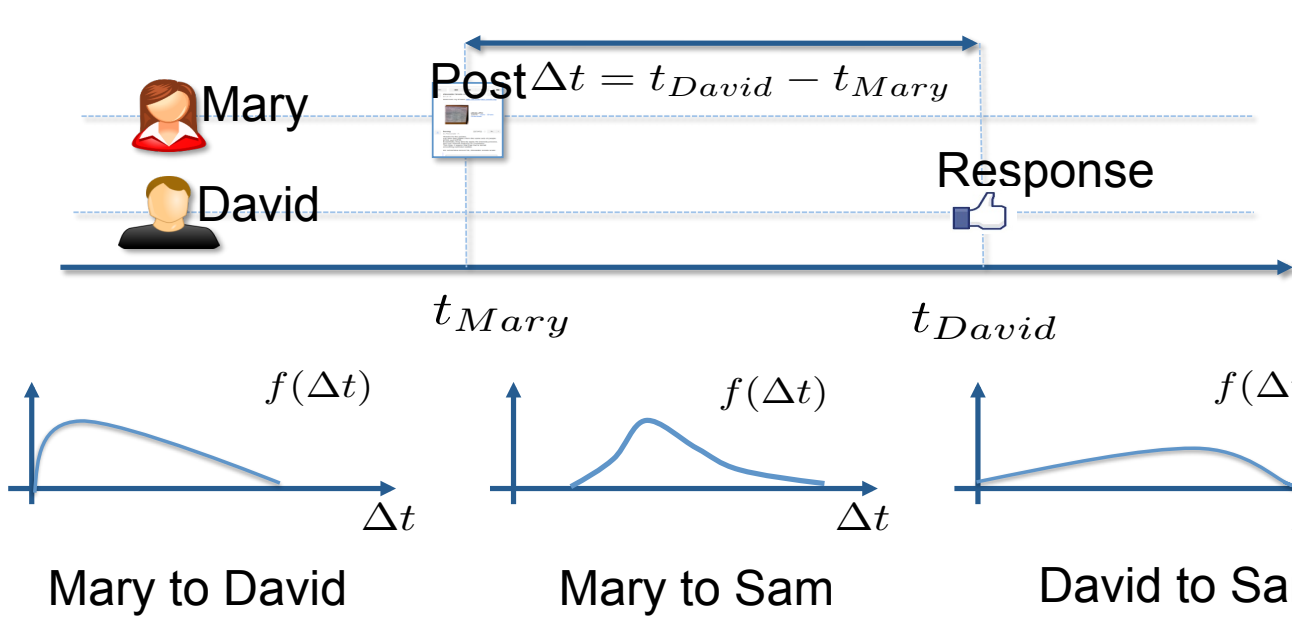


time-sensitive viral marketing



## CONTINUOUS-TIME INDEPENDENT CASCADE MODEL

- Infection : an event occurs to a node, e.g., adopting a product.
- Pairwise conditional density  $f_{ji}(t_j|t_i) = f_{ji}(t_i - t_j)$  over time



## ABSOLUTE INFECTION TIME PERSPECTIVE

- For each node  $i$ , let random variable  $t_i$  represent the infection time.
- The influence value of sources  $\mathcal{A}$  by time  $T$  is

$$\sigma(\mathcal{A}, T) = \mathbb{E} \left[ \sum_{i \in \mathcal{V}} \mathbb{I}\{t_i \leq T\} \right] = \sum_{i \in \mathcal{V}} \Pr\{t_i \leq T\}.$$

- Directed graphical model representation

$$p(\{t_i\}_{i \in \mathcal{V}}) = \prod_{i \in \mathcal{V}} p(t_i | \{t_j\}_{j \in \pi_i}), \quad \pi_i \text{ is the set of parents.}$$

- Marginalization

$$\Pr\{t_i \leq T\} = \int_0^\infty \cdots \int_{t=0}^T \cdots \int_0^\infty \left( \prod_{j \in \mathcal{V}} p(t_j | \{t_l\}_{l \in \pi_j}) \right) \left( \prod_{j \in \mathcal{V}} dt_j \right).$$

## INTER-EVENT TIME PERSPECTIVE

- Aim to calculate  $\mathbb{E} \left[ \sum_{i \in \mathcal{V}} \mathbb{I}\{t_i \leq T\} \right]$  directly.
- Mutually independent transmission times  $\tau_{ji} = t_i - t_j$ .

$$p(\{\tau_{ji}\}_{(j,i) \in \mathcal{E}}) = \prod_{(j,i) \in \mathcal{E}} f_{ji}(\tau_{ji}).$$

- A set of transmission times (or a particular configuration).

$$G := \{\tau_{ji}\}_{(j,i) \in \mathcal{E}} \sim p(\{\tau_{ji}\}_{(j,i) \in \mathcal{E}}).$$

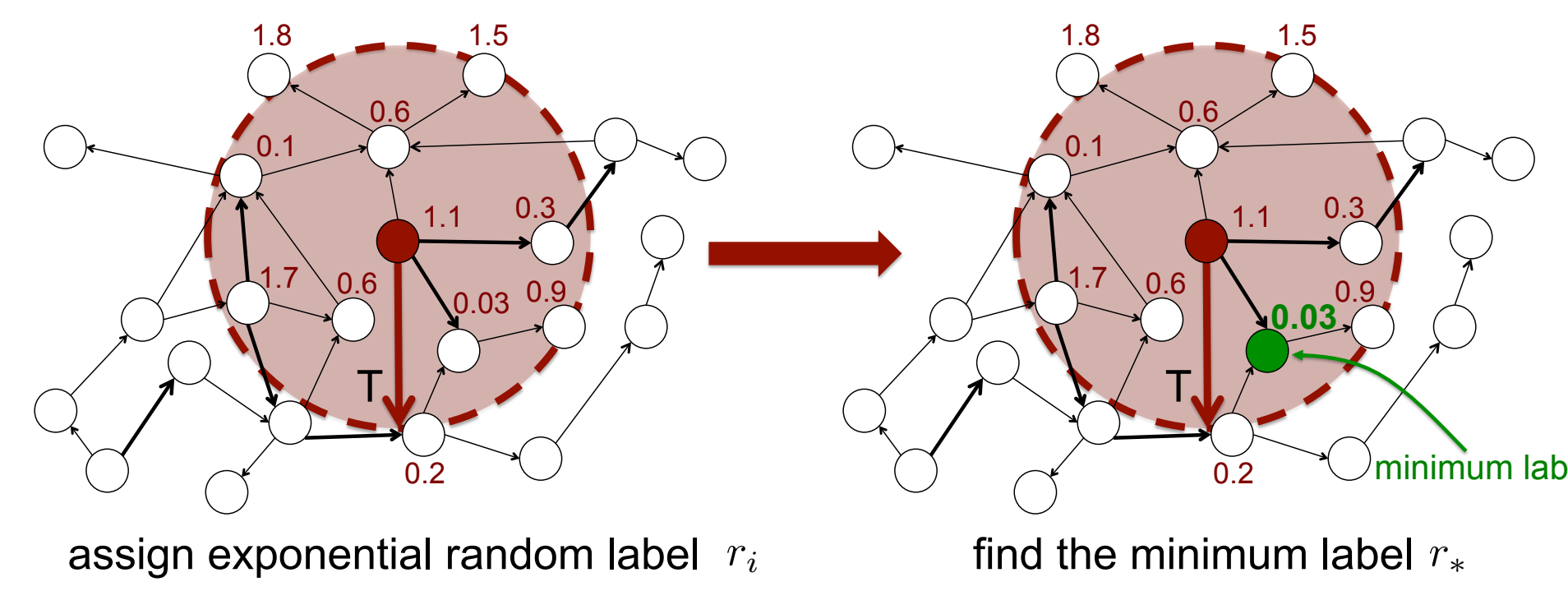
- Given  $G$ ,  $t_i$  is the length of shortest path from all sources in  $\mathcal{A}$  to  $i$ .
- Draw a set of  $G^l \sim p(\{\tau_{ji}\}_{(j,i) \in \mathcal{E}})$

$$\Pr\{t_i \leq T\} = \Pr\{\text{length of shortest path from } \mathcal{A} \rightarrow i \leq T\}.$$

- Naive simulation requires  $O(|\mathcal{V}|^2)$  time complexity.

## NEIGHBORHOOD SIZE ESTIMATION

Given  $n$  i.i.d random variable  $X^i \sim e^{-X}$ , the minimum  $X_* \sim ne^{-nx}$ .



- Find  $m$  such least labels,  $\{r_*^u\}_{u=1}^m$  to estimate  $|\mathcal{N}(\{j\}, T)| \approx \frac{m-1}{\sum_{u=1}^m r_*^u}$ , convert counting problem to estimation problem !

## NEIGHBORHOOD SIZE ESTIMATION

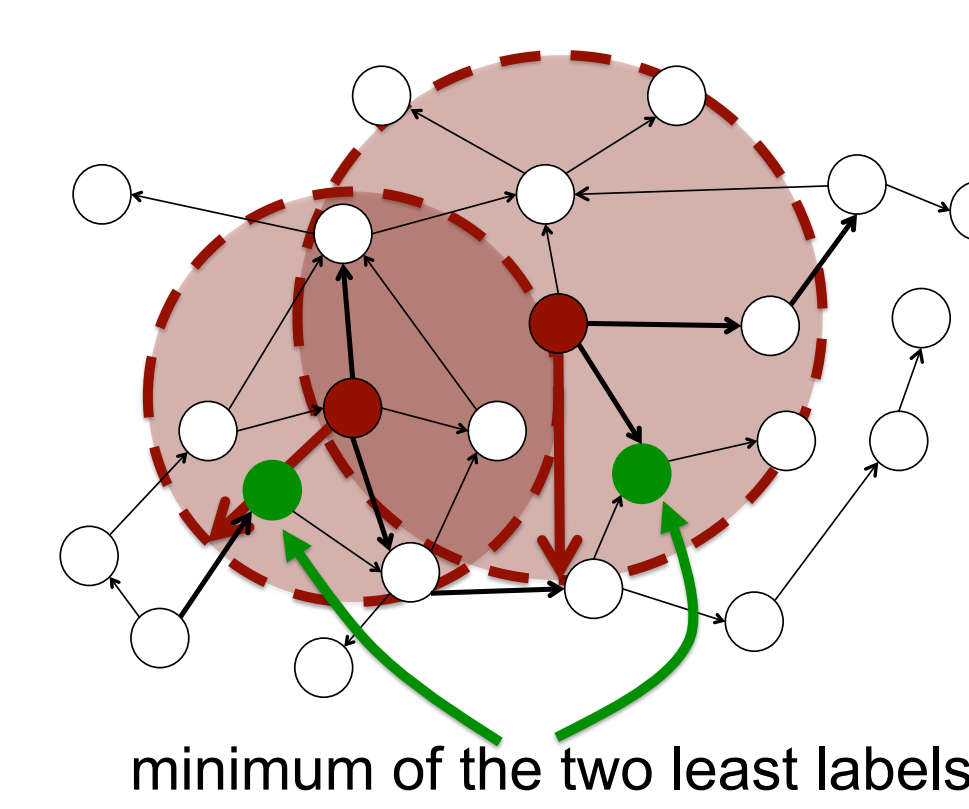
- Multiple sources  $\mathcal{A}$

$$\mathcal{N}(\mathcal{A}, T) = \bigcup_{s \in \mathcal{A}} \mathcal{N}(s, T).$$

- Reuse least-label list for each single source  $s \in \mathcal{A}$

$$r_* = \min_{i \in \mathcal{A}} \min_{j \in \mathcal{N}(i, T)} r_j$$

- Cohen's algorithm produces the lists for all nodes in time  $\tilde{O}(|\mathcal{E}|)$ .



## OVERALL ALGORITHM CONTINEST

Influence function

$$\sigma(\mathcal{A}, T) = \mathbb{E}_{\{\tau_{ij}\}_{(j,i) \in \mathcal{E}}} [|\mathcal{N}(\mathcal{A}, T)|] = \mathbb{E}_{\{\tau_{ij}\}} \mathbb{E}_{\{r^1, \dots, r^m\}} \left[ \frac{m-1}{\sum_{u=1}^m r_*^u} \right].$$

- Sample  $n$  sets of random transmission times  $\{\tau_{ij}^l\}_{(j,i) \in \mathcal{E}} \sim \prod_{(j,i) \in \mathcal{E}} f_{ji}(\tau_{ij}^l)$
- Given  $\{\tau_{ij}^l\}_{(j,i) \in \mathcal{E}}$ , sample  $m$  random labels  $\{r_i^u\}_{i \in \mathcal{V}} \sim \prod_{i \in \mathcal{V}} \exp(-r_i)$ .
- Estimate  $\sigma(\mathcal{A}, T)$  by  $\hat{\sigma}(\mathcal{A}, T) \approx \frac{1}{n} \sum_{l=1}^n ((m-1) / \sum_{u=1}^m r_*^u)$ .

## OVERALL ALGORITHM CONTINEST

**Theorem** : Draw the following number of sets of random transmission times

$$n \geq \frac{C\Lambda}{\epsilon^2} \log \left( \frac{2|\mathcal{V}|}{\delta} \right),$$

where  $\Lambda$  depends on  $\mathcal{A}$  and  $T$ , and for each set of random transmission times, draw  $m$  sets of random labels. Then  $|\hat{\sigma}(\mathcal{A}, T) - \sigma(\mathcal{A}, T)| \leq \epsilon$  uniformly for all  $\mathcal{A}$  with  $|\mathcal{A}| \leq C$ , with probability at least  $1 - \delta$ .

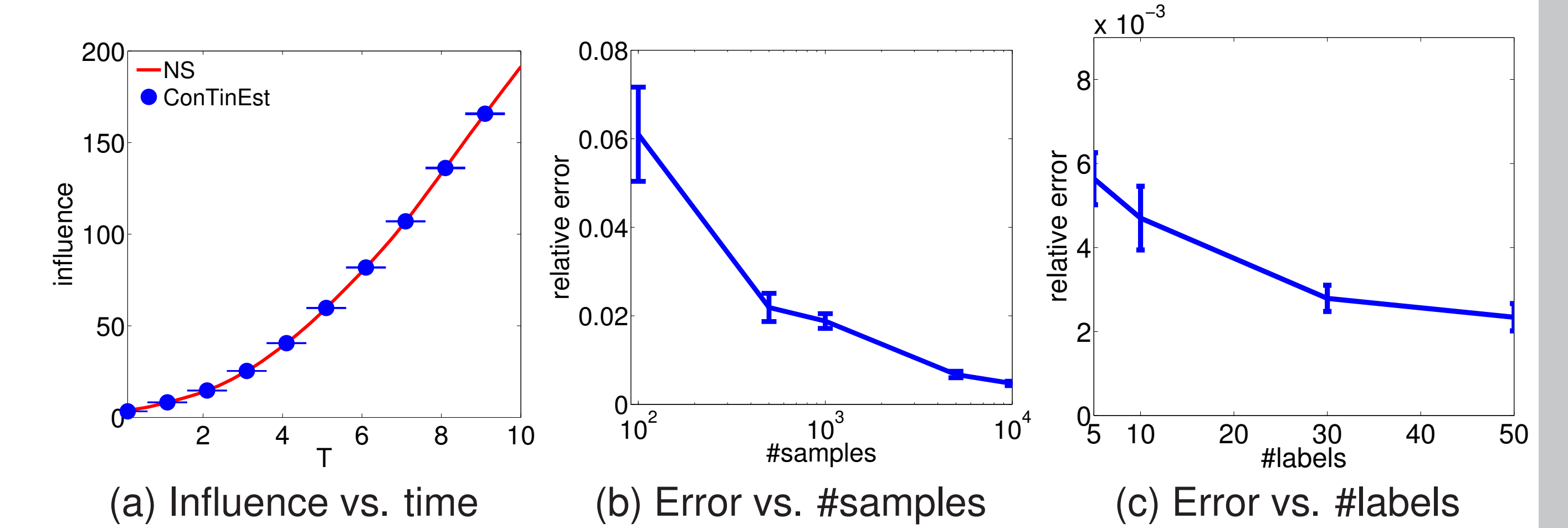
- Influence of larger source set  $\mathcal{A}$  at the longer time window  $T$  requires more samples in the worst case.
- In practice : small  $m = 5$  achieves good performance.

## INFLUENCE MAXIMIZATION CONTINEST

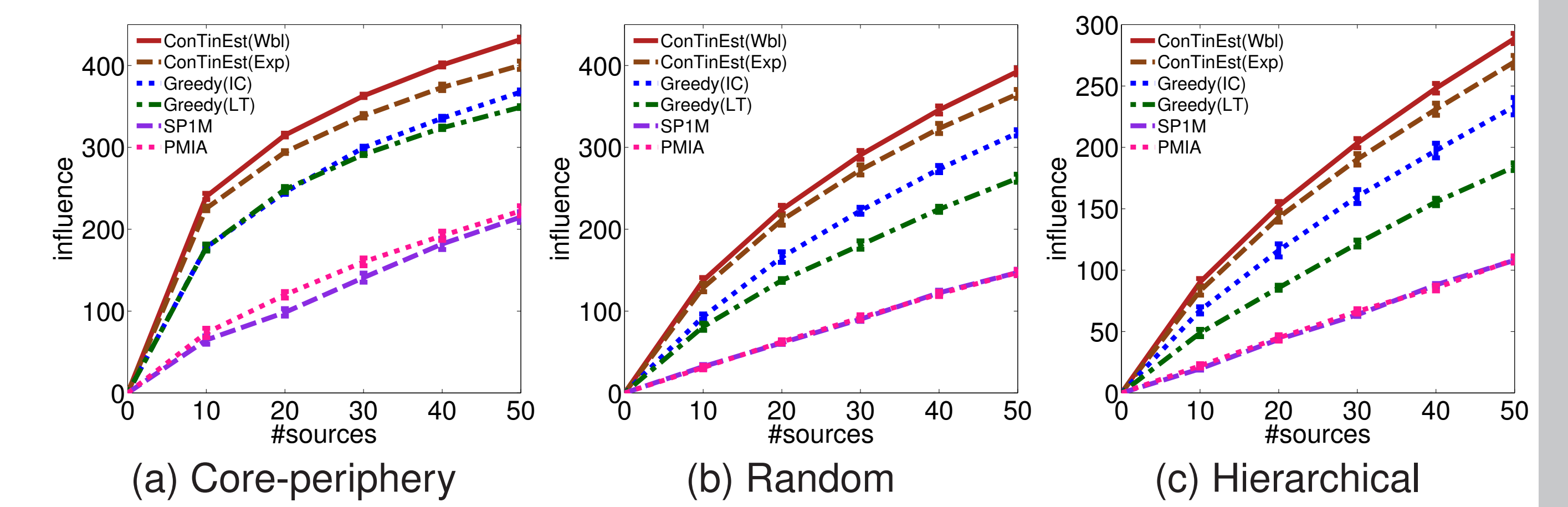
- Solve  $\mathcal{A}^* = \arg\max_{|\mathcal{A}| \leq C} \sigma(\mathcal{A}, T)$ , which is NP-hard in general.
- $\sigma(\mathcal{A}, T)$  is a non-negative, monotonic, submodular function.
- Greedy algorithm achieves  $(1 - 1/e)$  of the optimal value (OPT).
- The estimator  $\hat{\sigma}(\mathcal{A}, T)$  induces some error  $\epsilon$ .
- Theorem** : Suppose the influence  $\sigma(\mathcal{A}, T)$  for all  $\mathcal{A}$  with  $|\mathcal{A}| \leq C$  are estimated uniformly with error  $\epsilon$  and confidence  $1 - \delta$ , the greedy algorithm returns a set of sources  $\hat{\mathcal{A}}$  such that  $\sigma(\hat{\mathcal{A}}, T) \geq (1 - 1/e)OPT - 2C\epsilon$  with probability at least  $1 - \delta$ .

## EXPERIMENTAL EVALUATION : SYNTHETIC DATASET

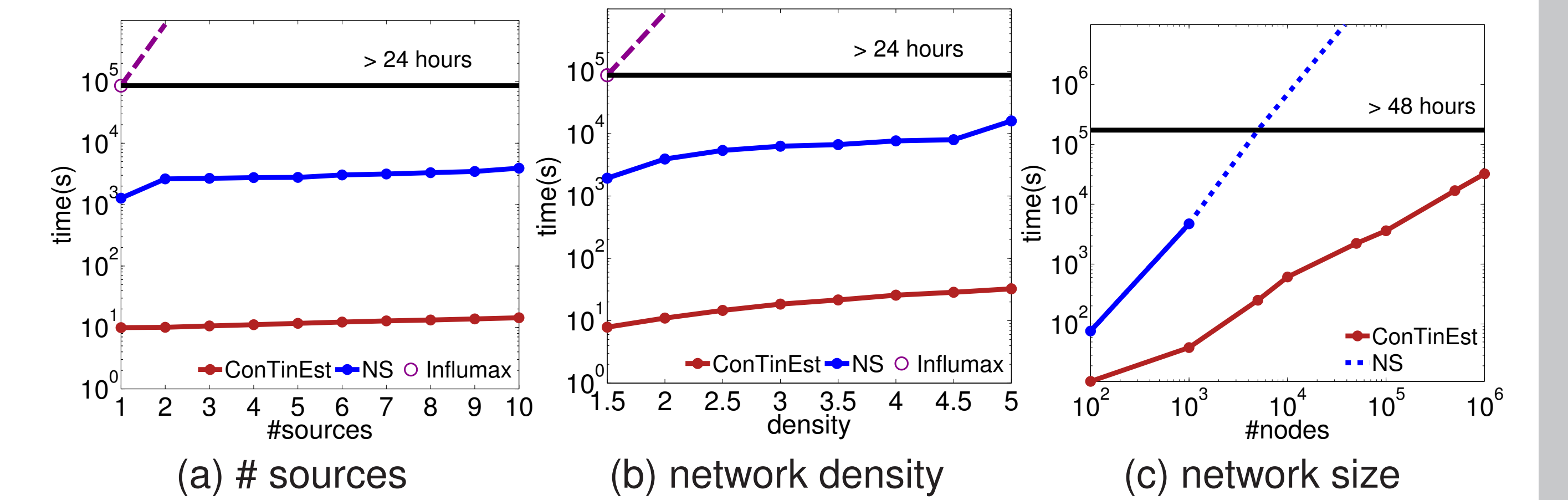
- Accuracy of the estimated influence (highest out-degree node).



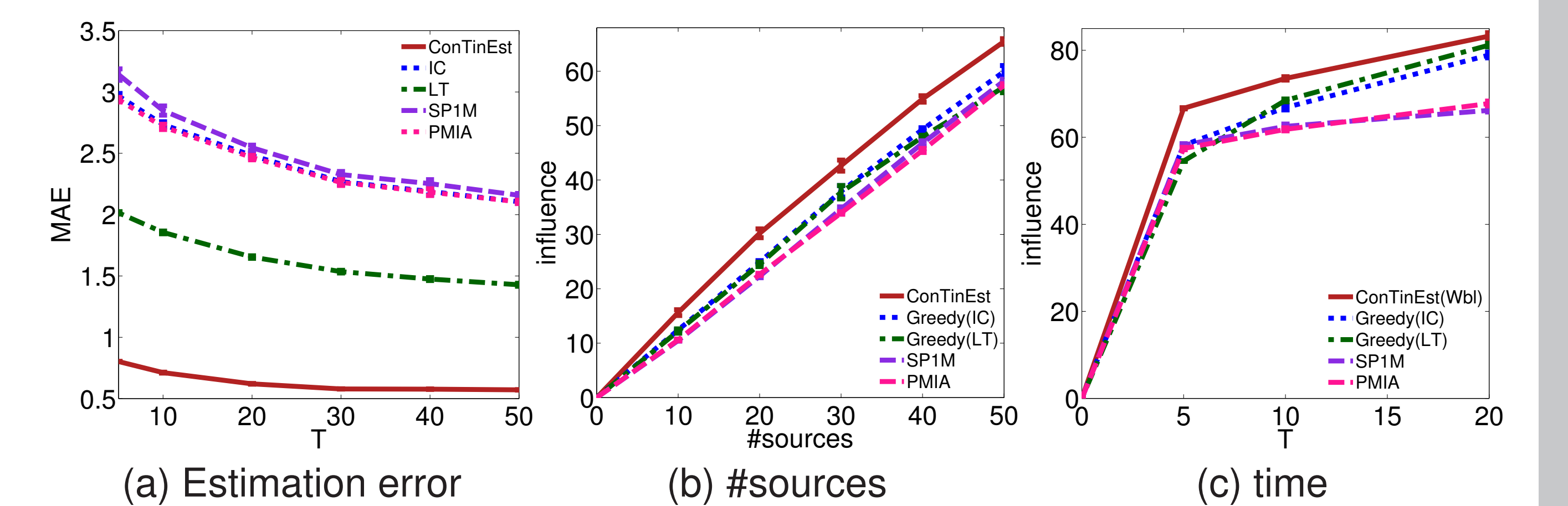
- Quality of the selected nodes for influence maximization.



- Scalability of influence maximization.



## EXPERIMENTAL EVALUATION : REAL DATASET



- 10,967 hyperlink cascades randomly splits into 80%-training and 20%-testing data.

- Infer network structures based on the training data.
- Evaluate estimated influence on the testing data.
  - Given node  $i$ ,  $\mathcal{C}(i)$  be the set of all cascades where  $i$  is the source.
  - Based on  $\mathcal{C}(i)$ , the total number of distinct nodes infected before  $T$  quantifies the real influence of node  $u$  up to time  $T$ .
  - Average across different cascades is the true influence.
- Repeat experiments for 10 times.