

Scaling Epidemic Inference on Contact Networks: Theory and Algorithms

Guanghui Min, Yinhan He, Chen Chen
University of Virginia



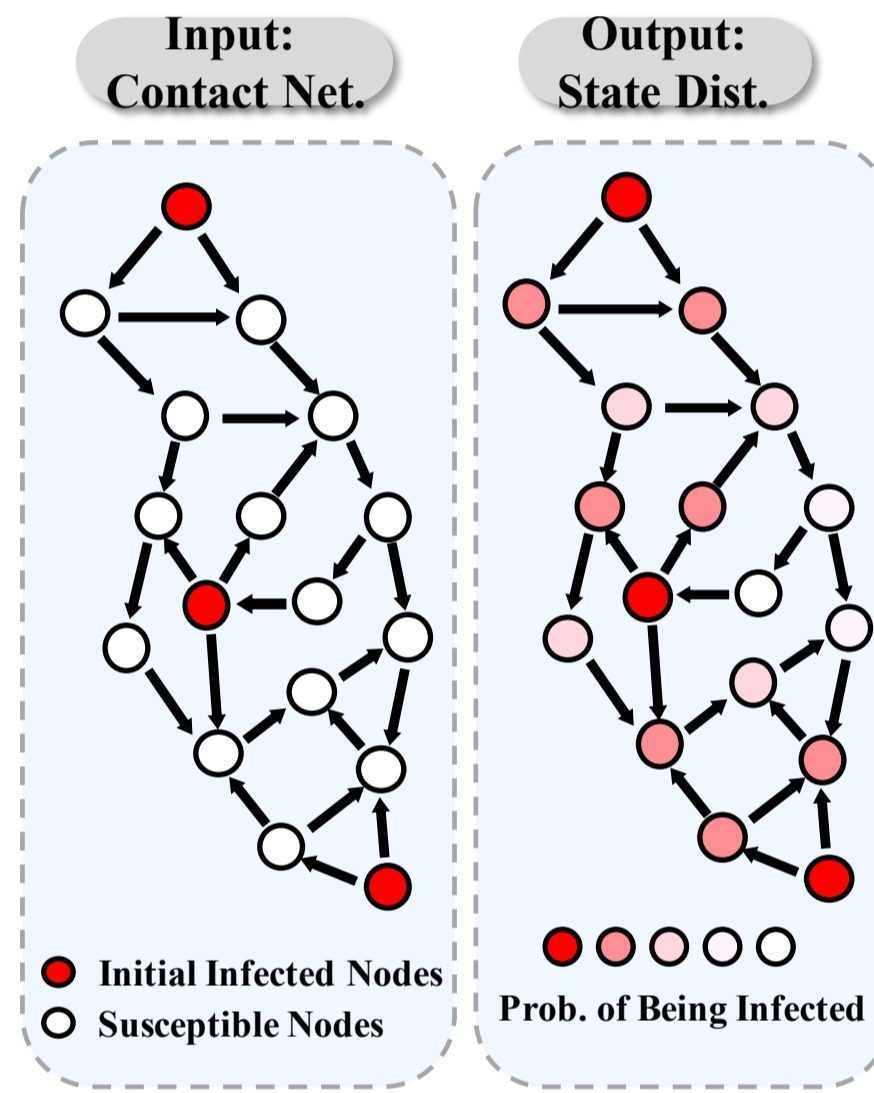
Motivations & Contributions

Challenges of Epidemic Inference

- High Computational Cost of Monte Carlo Simulation.** Current epidemic inference relies heavily on Monte Carlo (MC) simulation, which is statistically reliable but **computationally expensive**.
- Limitations of Population-Level Models.** Population-level models (e.g., SIR ODEs) assume homogeneous mixing and **cannot capture individual-level infection risk** or the impact of network topology.
- Shortcomings of Graph-Based Approaches.** Existing graph-based methods (like message passing or centrality) either **ignore infection dynamics** or lack theoretical guarantees on convergence and accuracy. Also the forward iteration is **computationally expensive**.

Our Contributions

- Theoretical Analysis of MC Variance.** Theoretically analyze how **network topology** and **epidemic parameters** affect MC variance;
- Development of the RAPID Algorithm.** Develop a **residual-aware propagation algorithm, RAPID** that matches MC accuracy at runtime comparable to a single MC simulation through asynchronous local updates;
- Comprehensive Empirical Validation.** Conduct experiments across six real-world networks, demonstrating robust **accuracy, scalability,** and **theoretical consistency**.



Theoretical Insight

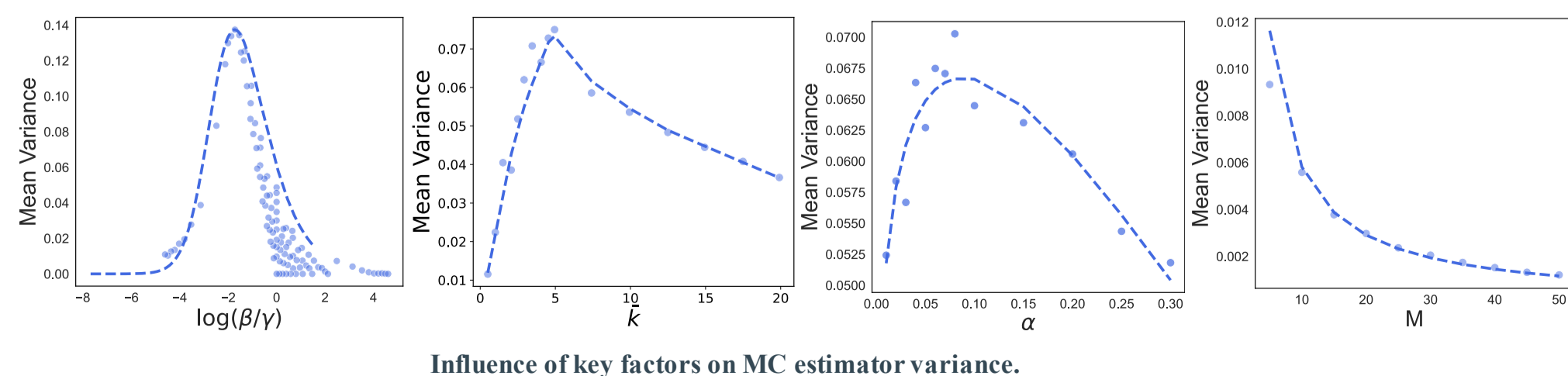
Theorem 3.1 quantifies how the variance of Monte Carlo (MC) estimators for node infection probability fundamentally depends on **epidemic parameters** (β, γ), **network structure** (average degree \bar{k} and diameter D), **initial infection fraction** α , and the **number of simulations** M .

It establishes a non-zero lower bound on the average estimator variance:

$$\frac{1}{N} \sum_{i=1}^N \text{Var}(\hat{p}_i - p_i) \gtrsim \frac{1}{2M} \min\{1 - (1 - p_0)^{c\bar{k}\alpha}, (1 - p_0)^{c\bar{k}\alpha}\},$$

where

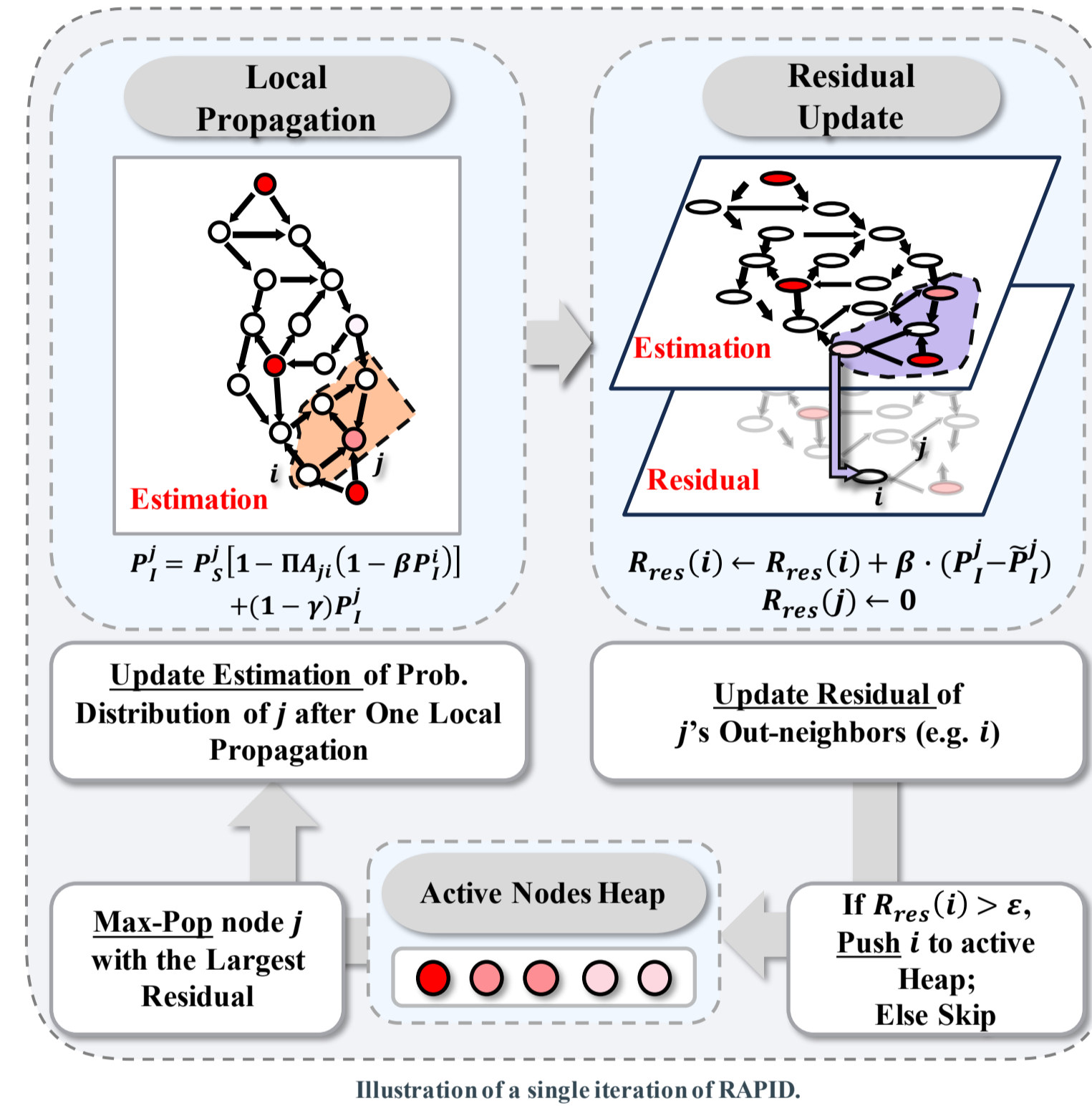
$$p_0 := \left(\frac{\beta}{\beta + \gamma}\right)^\ell, \quad \ell := \min\{D, \frac{\log N}{\log \bar{k}}\}$$



RAPID Algorithm

Core Idea

RAPID builds upon the Probabilistic Infection Dynamics (PID) **message-passing** equations and introduces a **residual-driven asynchronous propagation** mechanism that updates only where changes are significant.



Base: Message Passing Foundation

Each node i updates its infection probability P_i^i using local messages from its in-neighbors:

$$P_S^i(t+1) = P_S^i(t) \prod_{j \in \mathcal{V}} A_{ji} (1 - \beta P_I^j(t))$$

$$P_I^i(t+1) = P_S^i(t) [1 - \prod_{j \in \mathcal{V}} A_{ji} (1 - \beta P_I^j(t))] + (1 - \gamma) P_I^i(t)$$

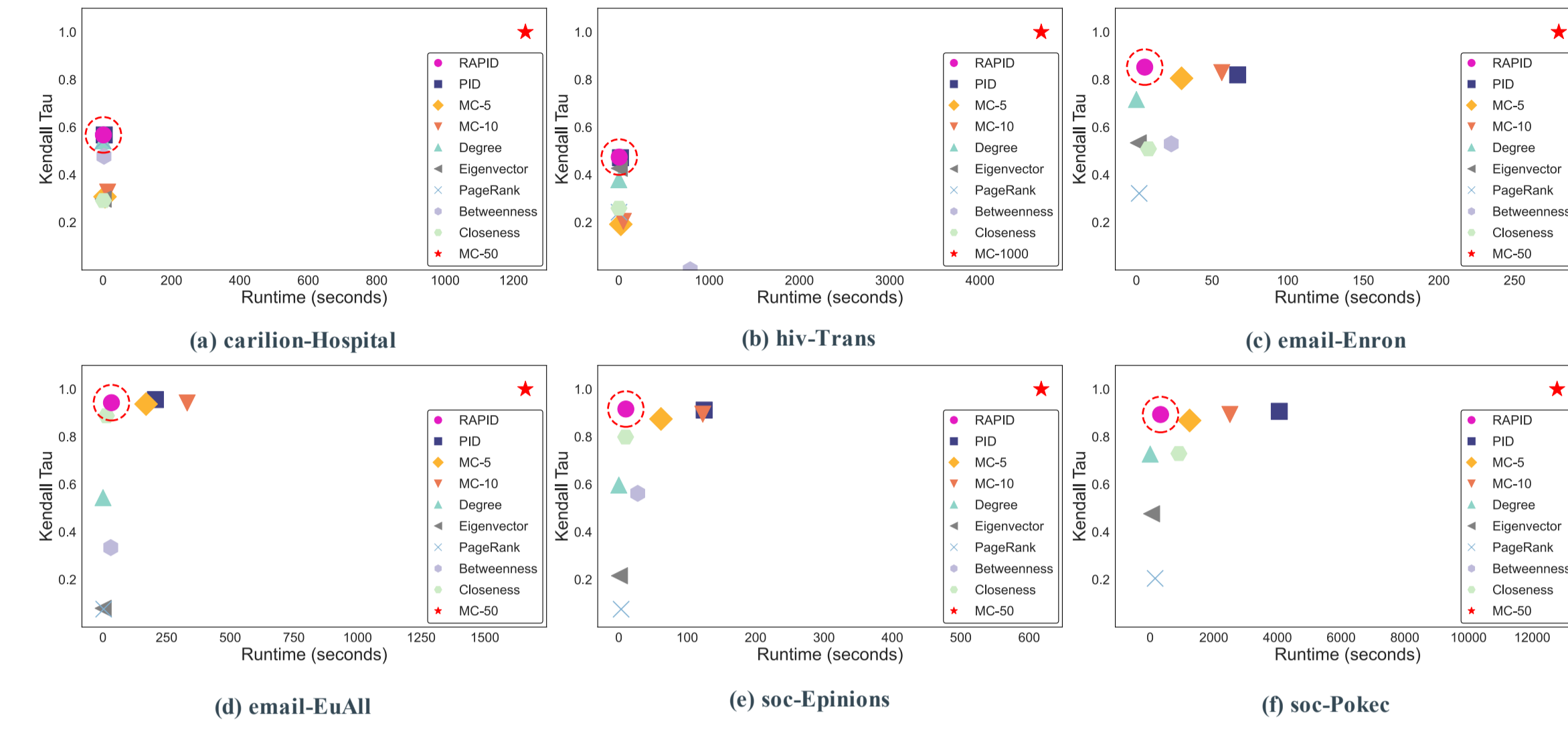
This standard PID update defines RAPID's computational base.

Residual-Driven Propagation

To quantify "how much information remains to be propagated", we define the propagation residual at node i :

$$R_{res}(i) = \beta \sum_{j \in \mathcal{V}} A_{ji} (P_I^j - \tilde{P}_I^j),$$

where \tilde{P}_I^j is the cached infection probability before the last update. This residual measures the unbalanced local influence, which is a first-order **Jacobian approximation** of the difference between current and previously propagated messages from the in-neighbors.



Trade-off between Kendall-Tau and Runtime across six datasets.

Experiments & Results

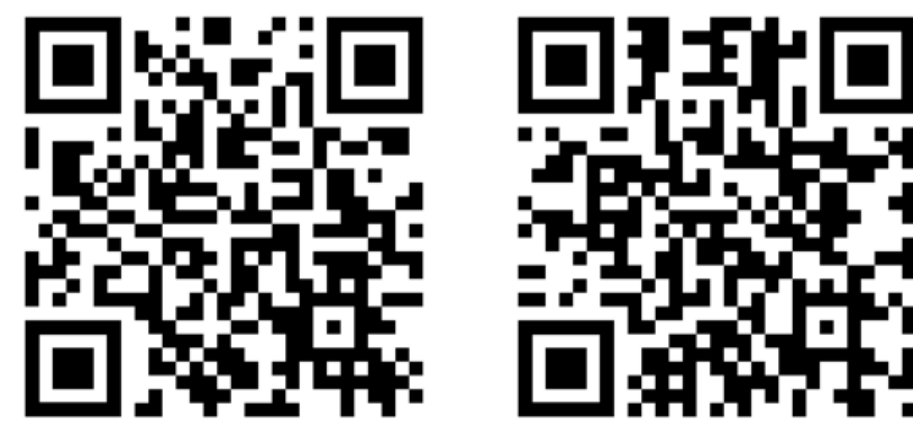
- Setup:** Six real-world directed networks (*carilion-Hospital*, *hiv-Trans*, *soc-Pokec*, etc.).
- Baselines:** MC-5/10/50, PID, and centrality heuristics.
- Metrics:** Kendall-tau Coefficient, Mean Absolute Error, Precision/Recall/F1, Runtime.
- Accuracy:** **RAPID** achieves the **lowest MAE** among baselines and **comparable Kendall-tau coefficient** to forward iteration baseline PID, benefiting from asynchronous updates that reduce the impact of small loops and accelerate convergence.
- Efficiency:** On six datasets, **RAPID** achieves an average speedup of 5.11 \times , 10.67 \times , and 8.52 \times over MC-5, MC-10, and PID, respectively. The speedup is more pronounced on denser graphs and its runtime remains comparable to a single-run Monte Carlo iteration across all datasets.
- Scalability:** Runtime grows sublinearly with network size, validating efficient inference.
- Robustness:** Performance remains stable across wide ranges of β , γ , and seed fractions α .

Takeaway

RAPID attains multi-run MC accuracy with consistent, dataset-wide runtime acceleration.

		carilion-Hospital ²	hiv-Trans ²	email-Enron	email-EuAll	soc-Epinions	soc-Pokec
MC-5	t	5.81 \pm 0.51	21.87 \pm 3.15	29.84 \pm 1.43	169.46 \pm 8.23	59.69 \pm 1.26	1241.00 \pm 18.79
	Δ	5.43 \times	5.16 \times	5.41 \times	5.06 \times	5.82 \times	3.78 \times
MC-10	t	13.46 \pm 1.14	49.94 \pm 1.74	56.45 \pm 0.79	330.97 \pm 6.81	122.31 \pm 2.44	2506.60 \pm 45.46
	Δ	12.58 \times	11.78 \times	10.24 \times	9.88 \times	11.91 \times	7.64 \times
MC-50	t	1234.73 \pm 13.13	4678.18 \pm 8.96	279.26 \pm 3.09	1659.57 \pm 23.55	614.58 \pm 2.36	12782.37 \pm 237.30
	Δ	1153.95 \times	1103.34 \times	50.66 \times	49.52 \times	59.86 \times	38.93 \times
PID	t	3.56 \pm 0.01	17.91 \pm 0.14	66.95 \pm 0.29	206.18 \pm 0.65	132.60 \pm 0.62	4056.89 \pm 6.40
	Δ	3.33 \times	4.22 \times	12.14 \times	6.15 \times	12.91 \times	12.36 \times
RAPID	t	1.07\pm0.00	4.24\pm0.03	5.51\pm0.04	33.50\pm0.05	10.27\pm0.09	328.28\pm0.66

Runtime comparison across datasets (seconds, lower is better). Δ indicates the speedup factor relative to RAPID, computed as $\Delta = \text{Baseline time} / \text{RAPID time}$. On *carilion-Hospital* and *hiv-Trans*, we adopt 1000-run MC simulations as the ground truth for acceptable estimator variance.



Paper

Code