

# Reconstructing Graph Diffusion History from a Single Snapshot

Ruizhong Qiu<sup>†</sup> Dingsu Wang<sup>†</sup> Lei Ying<sup>‡</sup> H. Vincent Poor<sup>§</sup> Yifang Zhang<sup>¶</sup> Hanghang Tong<sup>†</sup>

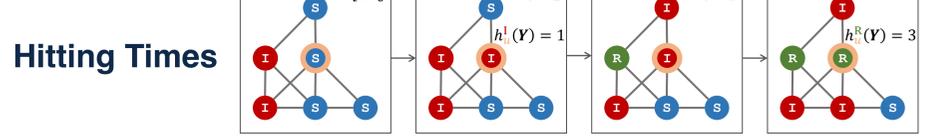
{rq5, dingsuw2, htong, zhang303}@illinois.edu leiying@umich.edu poor@princeton.edu



## HIGHLIGHTS

- **Problem definition:** Reconstructing Diffusion history from A Single Snapshot (DASH).
  - We do not assume knowing true diffusion parameters.
  - We do not assume having real histories as training data.
- **Theoretical insights:** Fundamental limitation of the MLE formulation.
  - Theorems 1 & 2  $\Rightarrow$  **Unavoidable** estimation error of diffusion parameters.
  - Theorem 3  $\Rightarrow$  The MLE formulation is **sensitive** to that estimation error.
- **Problem formulation:**
  - A novel **barycenter formulation** based on *hitting times*.
  - Provably **stable** against estimation error of diffusion parameters.
- **Proposed method:** Diffusion Hitting Times with Optimal proposal (DITTO).
  - Reducing the problem to estimating posterior expected hitting times via M-H MCMC;
  - Using a GNN to learn an **optimal** proposal to accelerate convergence of M-H MCMC.

## PROPOSED METHOD: DITTO



- First infection time:  $h_u^I(\mathbf{Y}) := \min\{T+1, \min\{t \geq 0: y_{t,u} \geq I\}\}$ .
- First recovery time:  $h_u^R(\mathbf{Y}) := \min\{T+1, \min\{t \geq 0: y_{t,u} \geq R\}\}$ .

## Stability of Posterior Expected Hitting Times

(Key theoretical observation)

- **Theorem 4.** Under SIR model and mild conditions, for any possible snapshot  $\mathbf{y}_T$ ,
 
$$\nabla_{\beta} \mathbb{E}_{Y \sim P_{\beta} | \mathbf{y}_T} [h_u^I(\mathbf{Y})] = O(1),$$

$$\nabla_{\beta} \mathbb{E}_{Y \sim P_{\beta} | \mathbf{y}_T} [h_u^R(\mathbf{Y})] = O(1).$$

➤ Stable even for small  $\beta$

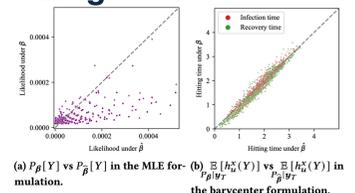


Figure 2: Sensitivity of the MLE formulation vs stability of the barycenter formulation.

- **Proof Idea:** Use Lemma 8 again to characterize  $P_{\beta}[\mathbf{Y}]$  and  $P_{\beta}[\mathbf{y}_T]$ .

## MLE Formulation $\rightarrow$ Barycenter Formulation

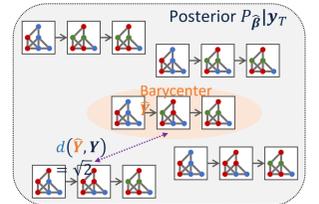
- Estimation error of  $\beta$  is **unavoidable**.
- The MLE formulation is **sensitive** to estimation error of  $\beta$ .
- Posterior expected *hitting times* are **stable** against estimation error of  $\beta$ .
- Our solution: A novel **barycenter formulation** based on *hitting times*.

- **History distance  $d$ :** (Euclidean distance with hitting times as coordinates)

$$d(\tilde{\mathbf{Y}}, \mathbf{Y}) := \sqrt{\sum_{u \in \mathcal{V}} \sum_{x \in \{I, R\}} (h_u^x(\tilde{\mathbf{Y}}) - h_u^x(\mathbf{Y}))^2}.$$

- **Barycenter formulation:** (Finding the **barycenter**  $\tilde{\mathbf{Y}}$  of the posterior distribution  $P_{\beta}[\mathbf{y}_T]$  w.r.t. the history distance  $d$ )

$$\min_{\tilde{\mathbf{Y}}} \mathbb{E}_{Y \sim P_{\beta} | \mathbf{y}_T} [d(\tilde{\mathbf{Y}}, \mathbf{Y})^2].$$



## Solution to the Barycenter Formulation

- Bias-variance decomposition:

$$\mathbb{E}_{Y \sim P_{\beta} | \mathbf{y}_T} [d(\tilde{\mathbf{Y}}, \mathbf{Y})^2] = \sum_{u \in \mathcal{V}} \sum_{x \in \{I, R\}} \left( (h_u^x(\tilde{\mathbf{Y}}) - \mathbb{E}_{Y \sim P_{\beta} | \mathbf{y}_T} [h_u^x(\mathbf{Y})])^2 + \text{Var}_{Y \sim P_{\beta} | \mathbf{y}_T} [h_u^x(\mathbf{Y})] \right).$$

- Variances are **constant** w.r.t.  $\tilde{\mathbf{Y}} \Rightarrow$  Optimal solution  $\tilde{\mathbf{Y}}$ :

$$h_u^x(\tilde{\mathbf{Y}}) = \text{round} \left( \mathbb{E}_{Y \sim P_{\beta} | \mathbf{y}_T} [h_u^x(\mathbf{Y})] \right), \quad x = I, R;$$

$$\hat{y}_{t,u} = \begin{cases} S, & \text{for } 0 \leq t < h_u^I(\tilde{\mathbf{Y}}); \\ I, & \text{for } h_u^I(\tilde{\mathbf{Y}}) \leq t < h_u^R(\tilde{\mathbf{Y}}); \\ R, & \text{for } h_u^R(\tilde{\mathbf{Y}}) \leq t \leq T. \end{cases}$$

## M-H MCMC for Posterior Expectation Estimation

- How to estimate  $\mathbb{E}_{Y \sim P_{\beta} | \mathbf{y}_T} [h_u^x(\mathbf{Y})]$ ?  $\triangleright$  Recall: **Intractable** to compute  $P_{\beta}[\mathbf{Y} | \mathbf{y}_T]$ .

- Our solution: **M-H MCMC** [1, 2].

1. Design a **proposal** distribution  $Q_{\theta}(\mathbf{y}_T)$  over possible histories.
2. Each step of M-H MCMC samples  $L$  histories  $\mathbf{Y}^{(s,i)} \sim Q_{\theta}(\mathbf{y}_T)$ ,  $i = 1, \dots, L$ .
3. Each previous history  $\mathbf{Y}^{(s-1,i)}$  is replaced by the new history  $\mathbf{Y}^{(s,i)}$  with probability:

$$\min \left\{ 1, \frac{P_{\beta}[\mathbf{Y}^{(s,i)} | \mathbf{y}_T] Q_{\theta}(\mathbf{y}_T) [\mathbf{Y}^{(s-1,i)}]}{P_{\beta}[\mathbf{Y}^{(s-1,i)} | \mathbf{y}_T] Q_{\theta}(\mathbf{y}_T) [\mathbf{Y}^{(s,i)}]} \right\} = \min \left\{ 1, \frac{P_{\beta}[\mathbf{Y}^{(s,i)}] Q_{\theta}(\mathbf{y}_T) [\mathbf{Y}^{(s-1,i)}]}{P_{\beta}[\mathbf{Y}^{(s-1,i)}] Q_{\theta}(\mathbf{y}_T) [\mathbf{Y}^{(s,i)}]} \right\} \triangleright \text{Tractable to compute}$$

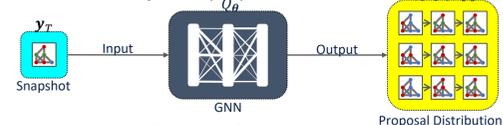
- The Markov chain  $\{\mathbf{Y}^{(s,i)}\}$  **provably converges** to the posterior distribution  $P_{\beta}[\mathbf{y}_T]$ .

$$\mathbb{E}_{Y \sim P_{\beta} | \mathbf{y}_T} [h_u^x(\mathbf{Y})] \approx \frac{1}{L} \sum_{i=1}^L h_u^x(\mathbf{Y}^{(s,i)}), \quad s \rightarrow +\infty.$$

## Learning an Optimal Proposal for M-H MCMC

- The **convergence rate** of M-H MCMC depends critically on the proposal  $Q_{\theta}$ .
- $\triangleright Q_{\theta}(\mathbf{y}_T)$  closer to  $P_{\beta}[\mathbf{y}_T] \Rightarrow$  Higher rate of convergence.

- Our solution: Use a GNN to learn an **optimal** proposal.

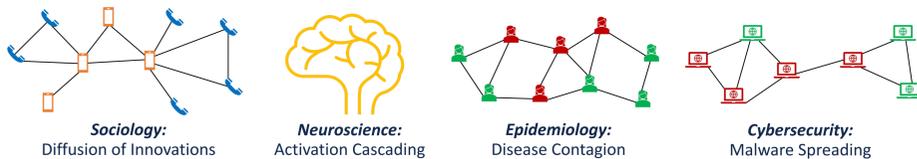


- Objective function: (a corollary of our Theorem 5; see our paper for details)

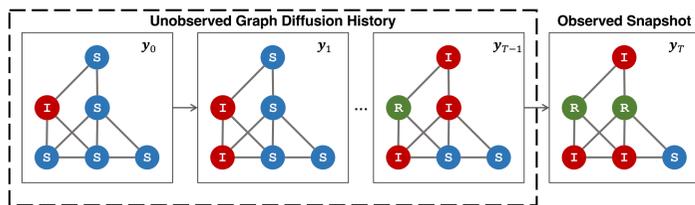
$$\min_{\theta} \mathbb{E}_{Y \sim P_{\beta}} [-\log Q_{\theta}(\mathbf{y}_T) | \mathbf{Y}].$$

## INTRODUCTION

### Diffusion on Graphs



### Problem Definition



$\mathcal{X} = \{S \text{ Susceptible}, I \text{ Infected}, R \text{ Recovered}\}$

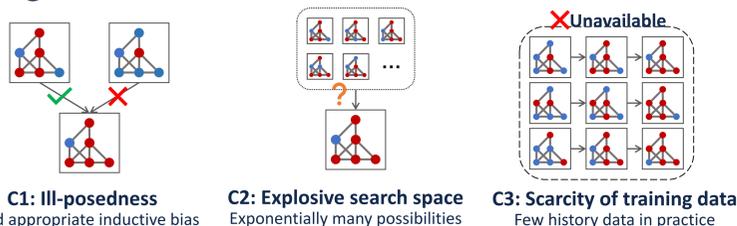
**Problem (DASH):** Reconstructing Diffusion history from A single Snapshot.

**Input:** (i) graph  $(\mathcal{V}, \mathcal{E})$ ; (ii) timespan  $T$  of interest; (iii) final snapshot  $\mathbf{y}_T \in \mathcal{X}^{\mathcal{V}}$ ; (iv) initial distribution  $P[\mathbf{y}_0]$ .

**Output:** reconstructed complete diffusion history  $\tilde{\mathbf{Y}} = [\tilde{\mathbf{y}}_0, \dots, \tilde{\mathbf{y}}_{T-1}, \mathbf{y}_T]^T \in \mathcal{X}^{T \times \mathcal{V}}$ .

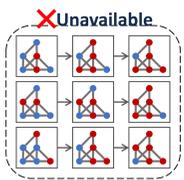
$\triangleright$  We do not assume knowing true diffusion parameters.

### Challenges of the DASH Problem

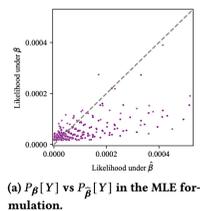


### Previous Methods & Their Limitations

Supervised time series imputation is impractical due to the scarcity of training data.



Maximum likelihood estimation (MLE) is sensitive to estimation error of diffusion parameters (our Theorems 1 & 2).



(a)  $P_{\beta}[\mathbf{Y}]$  vs  $P_{\beta}[\mathbf{y}_T]$  in the MLE formulation.

## REVISITING DIFFUSION HISTORY MLE

### NP-Hardness of Diffusion Parameter Estimation

- To estimate diffusion parameters  $\beta$ , a conventional approach is MLE:  $\max_{\beta} P_{\beta}[\mathbf{y}_T]$ . (\*)

- **Theorem 1 (informal):** Computing the probability  $P_{\beta}[\mathbf{y}_T]$  is NP-hard.  $\triangleright O\left(\binom{T+1}{2}^{n+m}\right)$  time.

- Think deeper: Is there an algo for  $\hat{\beta}$  MLE without computing  $P_{\beta}[\mathbf{y}_T]$ ?

- **Theorem 2 (informal):** Diffusion parameter MLE (\*) is NP-hard.

$\Rightarrow$  **Implication:** Estimation error of  $\beta$  is unavoidable.

### Sensitivity to Estimation Error of Diffusion Parameters

- **MLE formulation** for diffusion history reconstruction:

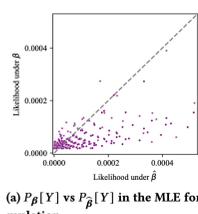
$$\max_{\beta \in \text{supp}(P | \mathbf{y}_T)} P_{\beta}[\tilde{\mathbf{Y}}].$$

- **Theorem 3.** Under the SIR model and mild conditions, for all possible history  $\mathbf{Y}$ , we have:

$$\frac{\partial}{\partial \beta^I} P_{\beta}[\mathbf{Y}] = \theta \left( \frac{1}{\beta^I} \right) P_{\beta}[\mathbf{Y}], \quad \frac{\partial}{\partial \beta^R} P_{\beta}[\mathbf{Y}] = \theta \left( \frac{1}{\beta^R} \right) P_{\beta}[\mathbf{Y}].$$

- Real-world diffusion typically has small true  $\beta$ .

$\Rightarrow$  **Implication:** MLE formulation is sensitive to estimation error of  $\hat{\beta}$ .



(a)  $P_{\beta}[\mathbf{Y}]$  vs  $P_{\beta}[\mathbf{y}_T]$  in the MLE formulation.

## MAIN EXPERIMENTS

### Performance for Real-World Diffusion

Table 4: Results for real-world diffusion. "OOM" indicates "out of memory."

Type	Method	BrFarmers		Pol		Covid		Hebrew	
		F1T	NRMSSE <sub>L</sub>	F1T	NRMSSE <sub>L</sub>	F1T	NRMSSE <sub>L</sub>	F1T	NRMSSE <sub>L</sub>
Supervised (w/ estimated $\hat{\beta}$ )	GCN	.5409	.6660	.4458	.4946	.3162	.5214	.3350	.6070
	GIN	.4548	.6565	.5203	.4767	.3226	.4951	.3704	.7816
	BRTS	.5207	.3995	OOM	OOM	.3524	.3333	.3120	.6584
	GRIN	.8003	.2425	.6518	.3731	.5448	.3040	.5916	.2212
	SPIN	.8268	.2084	OOM	OOM	.5917	.2932	.5178	.3330
MLE	DHREC	.6131	.4150	.7023	.3398	.3540	.6023	.6251	.4169
	CRI	.6058	.4444	.7468	.2942	.4170	.5487	.5344	.3552
Barycenter	DITTO (ours)	.8206	.2142	.7471	.2903	.6240	.2637	.6411	.2983

$\triangleright$  BrFarmers is very close to SI.

**DITTO:** Consistently strong performance across all datasets.

MLE/Supervised: Bad when real diffusion deviates from SI/SIR.

### Comparison with MLE-Based Methods

Table 5: Comparison with MLE-based methods on synthetic SI and SIR diffusion. \*We use GRIN trained with true  $\beta$  as the ideal performance and calculate Gap w.r.t. this ideal performance.

Type	Method	BA-SI			ER-SI			Oregon2-SI			Prot-SI				
		F1T	Gap <sub>L</sub>	NRMSSE <sub>L</sub>	F1T	Gap <sub>L</sub>	NRMSSE <sub>L</sub>	F1T	Gap <sub>L</sub>	NRMSSE <sub>L</sub>	F1T	Gap <sub>L</sub>	NRMSSE <sub>L</sub>		
Ideal	GCN	.7867	-.1692	.7827	.7827	-.1692	.7827	.7827	-.1692	.7827	.7827	-.1692	.7827		
	DITTO	.7867	-.1692	.7827	.7827	-.1692	.7827	.7827	-.1692	.7827	.7827	-.1692	.7827		
MLE	DHREC	.5900	35.43%	.4722	179.68%	.5500	27.88%	.4423	78.06%	.4044	24.85%	.4475	171.23%	.6268	22.30%
	CRI	.5994	21.81%	.3556	98.35%	.4129	19.63%	.3109	25.16%	.5761	28.20%	.3576	116.60%	.5738	28.87%
Barycenter	DITTO (ours)	.7783	1.07%	.1633	-3.49%	.7734	-1.42%	.1679	-32.41%	.7928	1.20%	.1707	3.39%	.7929	1.71%

**DITTO:** Stably achieves the strongest performance.

MLE: Performance vary largely across datasets due to sensitivity.

## ACKNOWLEDGEMENTS

