

## SPREADING PROCESSES ARE UBIQUITOUS

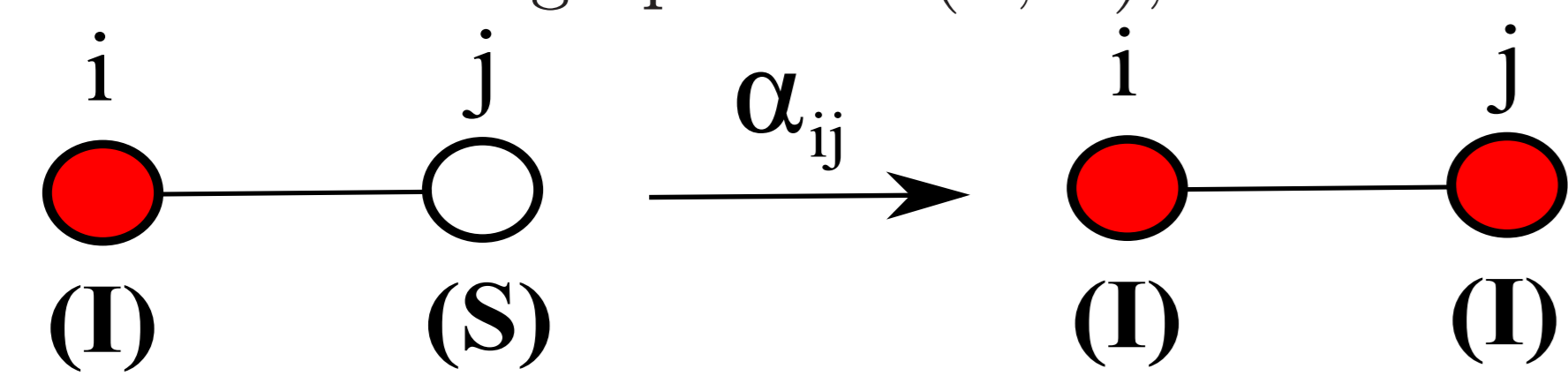
Epidemic spreading    Infrastructure failures    Financial contagion    Neural cascades    Malware propagation



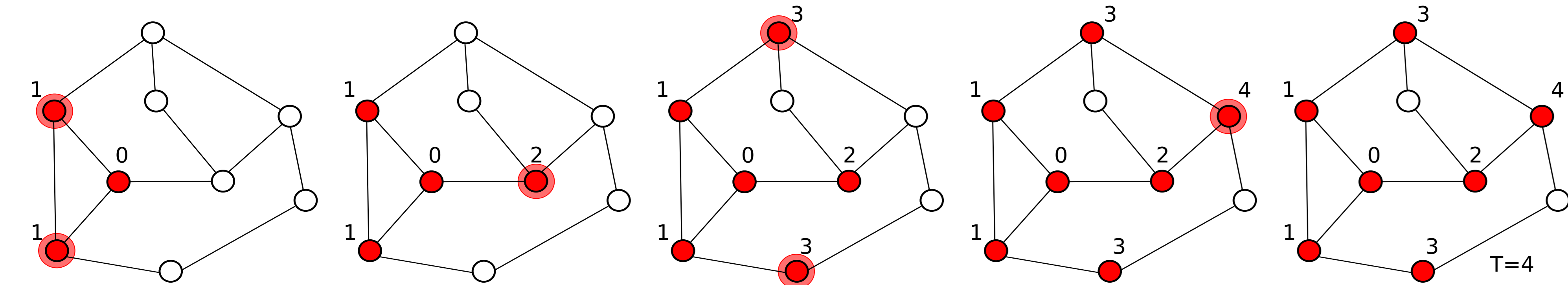
Common feature: **unknown models and parameters**

## EXAMPLE: SUSCEPTIBLE-INFECTED MODEL

The **SI model** is slightly more general than the popular Independent Cascade (IC) model – what follows applies to IC models as well. On a graph  $G = (V, E)$ , at each discrete time interval  $t \rightarrow t + 1$ :

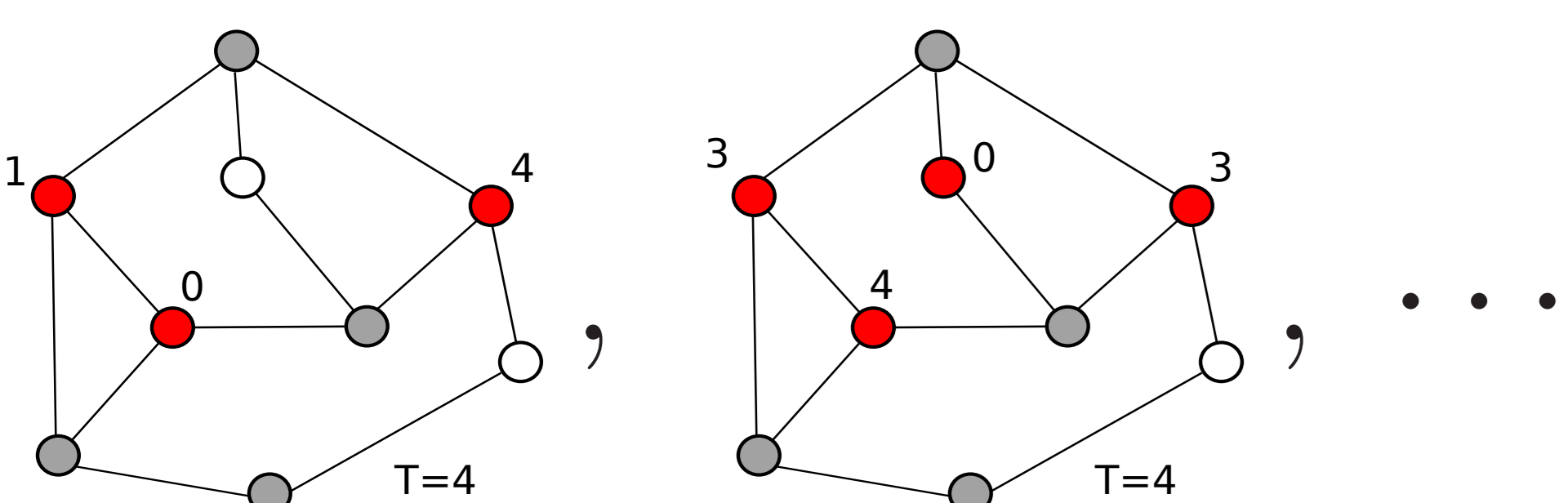


Denote  $\{\alpha_{ij}\}_{(ij) \in E} \equiv G_\alpha$ . **Cascade**  $\Sigma^c$ : collection of activation times  $\{\tau_i^c\}_{i \in V}$  until time horizon  $T$



## PROBLEM FORMULATION

Each cascade is divided into **observed** ( $\mathcal{O}$ ) and **hidden** ( $\mathcal{H}$ ) parts,  $\Sigma^c = \Sigma_{\mathcal{O}}^c \cup \Sigma_{\mathcal{H}}^c$ . **Examples**:  $|\mathcal{H}| = H$  hidden nodes not reporting their activation times, or snapshots of the network at a subset of times.



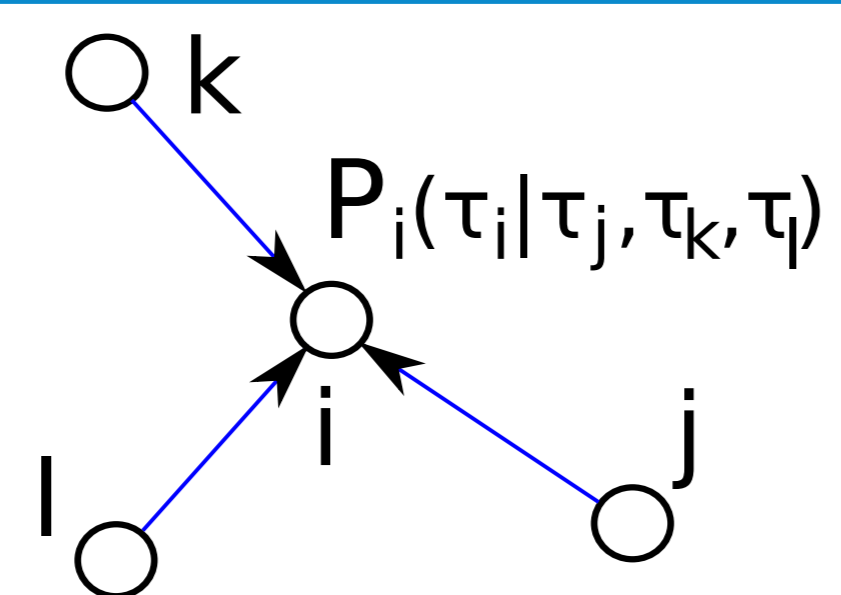
**Formulation of the problem:**

Given  $M$  **partially observed** cascades  $\Sigma_{\mathcal{O}} = \cup_{c=1}^M \Sigma_{\mathcal{O}}^c$ , **reconstruct model parameters**  $\{\alpha_{ij}^*\}_{(ij) \in E} \equiv G_{\alpha^*}$  used to generate data

## SPECIAL CASE: FULL OBSERVATIONS

**Full observations**  $\Sigma_{\mathcal{O}} = \Sigma$ : **easy**

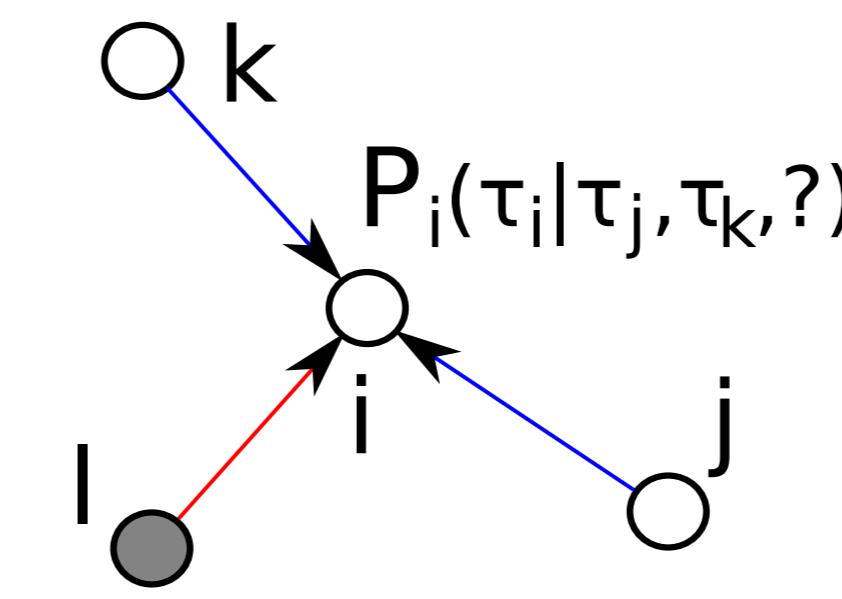
Maximize the likelihood function  $P(\Sigma | G_\alpha) = \prod_{i \in V} \prod_{1 \leq c \leq M} P_i(\tau_i^c | \Sigma^c, G_\alpha)$   
 $\hat{G}_{\alpha^*} = \arg \min (-\ln P(\Sigma | G_\alpha))$ : local convex optimization for each  $i \in V$



## MAXIMIZATION OF THE LIKELIHOOD: INTRACTABLE

**Partial observations**  $\Sigma_{\mathcal{O}} \neq \Sigma$ : **hard**

Maximization of the likelihood marginalized over unknown information:  
 $P(\Sigma_{\mathcal{O}} | G_\alpha) = \sum_{\{\tau_h^c\}_{h \in \mathcal{H}}} P(\Sigma | G_\alpha)$ , computational complexity  $\propto T^H$



## Proposed speed-up: Heuristic Two-Stage Algorithm (HTS)

1. Complete missing  $\{\tau_h^c\}_{h \in \mathcal{H}}$  by most probable values  $\hat{\Sigma}_{\mathcal{H}} = \arg \max P(\Sigma | \hat{G}_\alpha)$  using MC sampling.
2. Solve the “full observations” problem using  $\Sigma = \Sigma_{\mathcal{O}} \cup \hat{\Sigma}_{\mathcal{H}}$ .
3. Iterate steps 1 and 2 until global convergence of the algorithm.  $\square$  (Still very slow.)

## SOLUTION: DYNAMIC MESSAGE-PASSING ALGORITHM

**Dynamic Message-Passing (DMP)** equations allow to approximately solve the dynamics, i.e. compute the marginal probability  $m^i(t)$  of activation of node  $i$  at time  $t$  for each  $i \in V$  in time  $O(|E|T)$

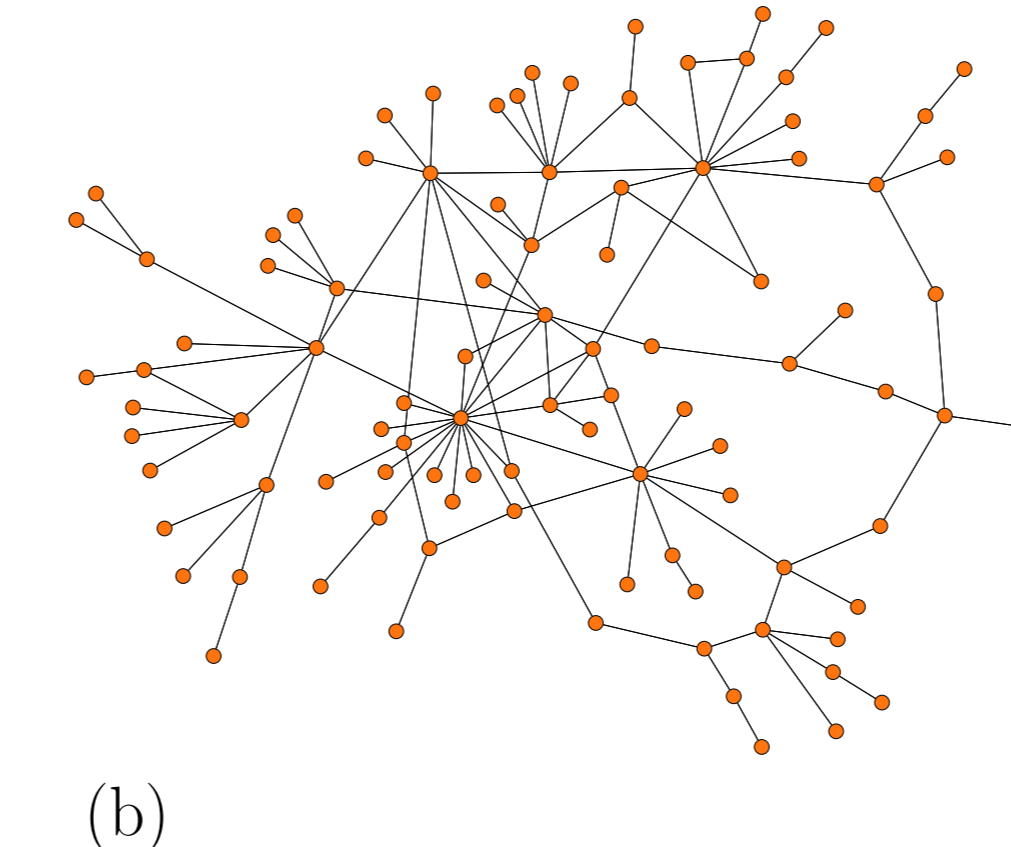
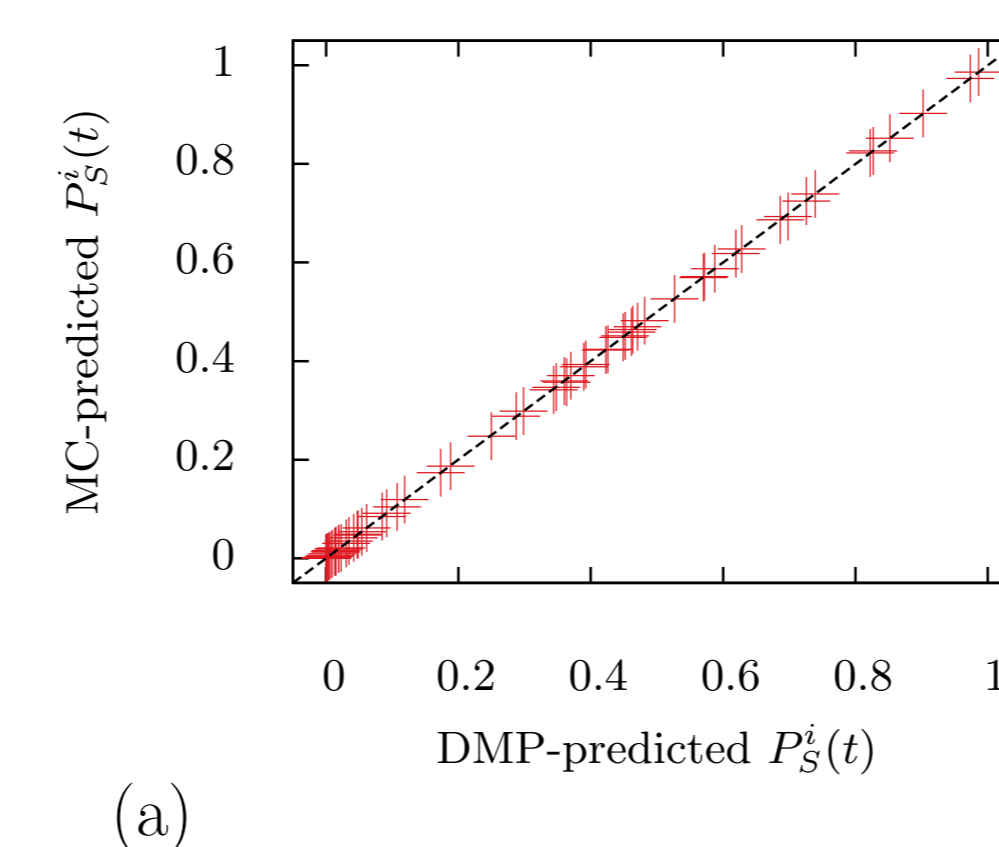
$$m^i(t) = P_S^i(t-1) - P_S^i(t),$$

$$P_S^i(t) = P_S^i(0) \prod_{k \in \partial i} \theta^{k \rightarrow i}(t),$$

$$\theta^{k \rightarrow i}(t) = \theta^{k \rightarrow i}(t-1) - \alpha_{ki} \phi^{k \rightarrow i}(t-1),$$

$$\phi^{k \rightarrow i}(t) = (1 - \alpha_{ki}) \phi^{k \rightarrow i}(t-1)$$

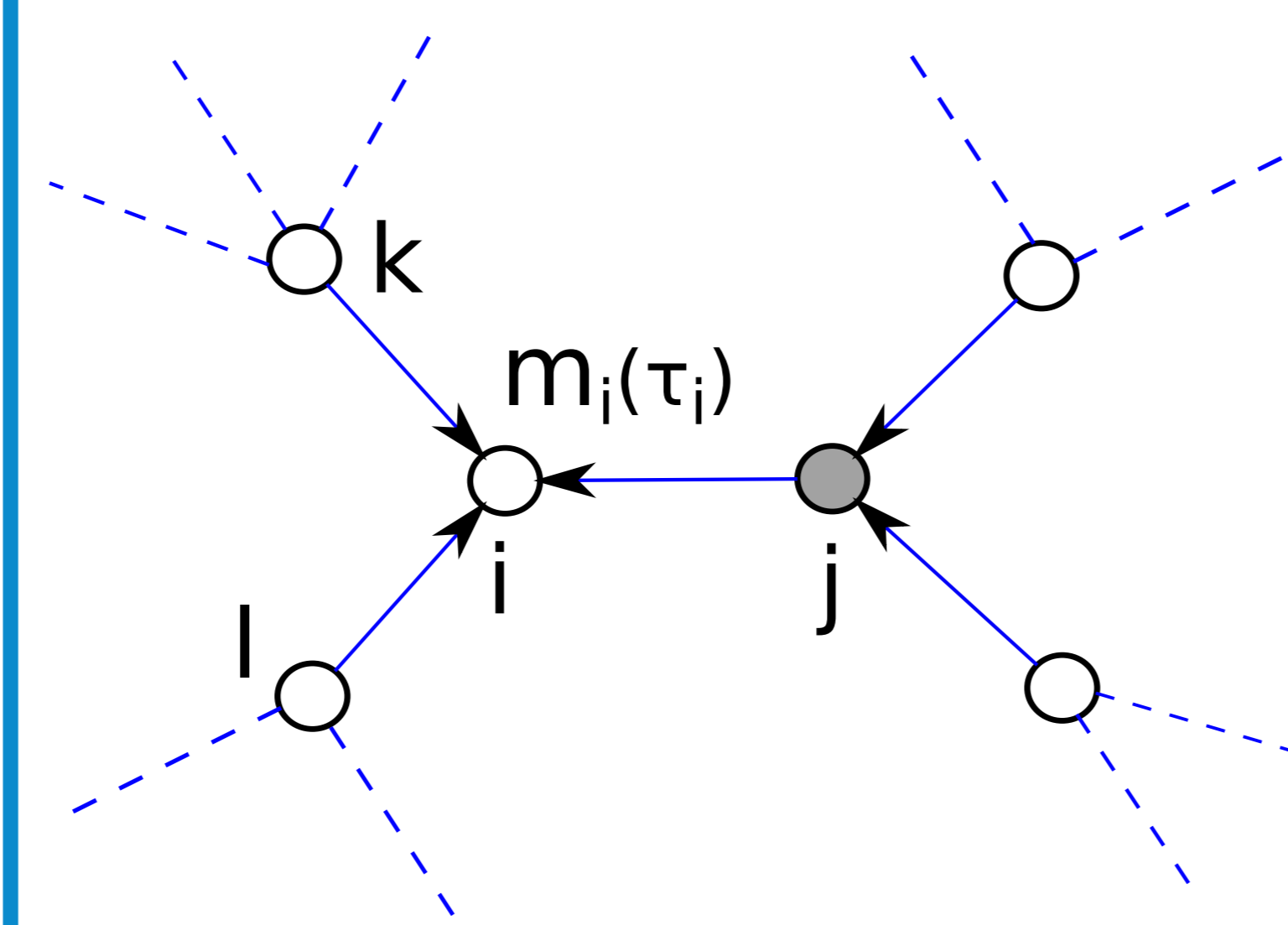
$$+ P_S^k(0) \prod_{l \in \partial k \setminus i} \theta^{l \rightarrow k}(t-1) - P_S^k(0) \prod_{l \in \partial k \setminus i} \theta^{l \rightarrow k}(t)$$



**Theorem:** The quantities  $P_S^i(t)$  are exact on tree graphs, and give lower bounds on values of marginal probabilities for general loopy graphs.  $\square$  (Accurate in practice.)

**Key idea:** approximation of the likelihood with marginal probabilities

$$P(\Sigma_{\mathcal{O}} | G_\alpha) \approx \prod_{c=1}^M \prod_{i \in \mathcal{O}} [m^i(\tau_i^c | G_\alpha) \mathbb{1}_{\tau_i^c \leq T} + P_S^i(\tau_i^c | G_\alpha) \mathbb{1}_{\tau_i^c = T}]$$



Each  $m^i(\tau_i^c)$  summarizes the effect of all possible propagation paths. Minimization of the “free energy”  $f_{\text{DMP}} = -\ln P(\Sigma_{\mathcal{O}} | G_\alpha) = \sum_{i \in \mathcal{O}} f_{\text{DMP}}^i$  will yield the **most likely consensus** among the ensemble of parameters.

Gradient  $\partial f_{\text{DMP}}^i / \partial \alpha_{rs}$  through the DMP eqs. for derivatives  $p_{rs}^{k \rightarrow i}(t) \equiv \partial \theta^{k \rightarrow i}(t) / \partial \alpha_{rs}$ , e.g.

$$\frac{\partial P_S^i(t)}{\partial \alpha_{rs}} = P_S^i(0) \sum_{k \in \partial i} p_{rs}^{k \rightarrow i}(t) \prod_{l \in \partial i \setminus k} \theta^{l \rightarrow i}(t)$$

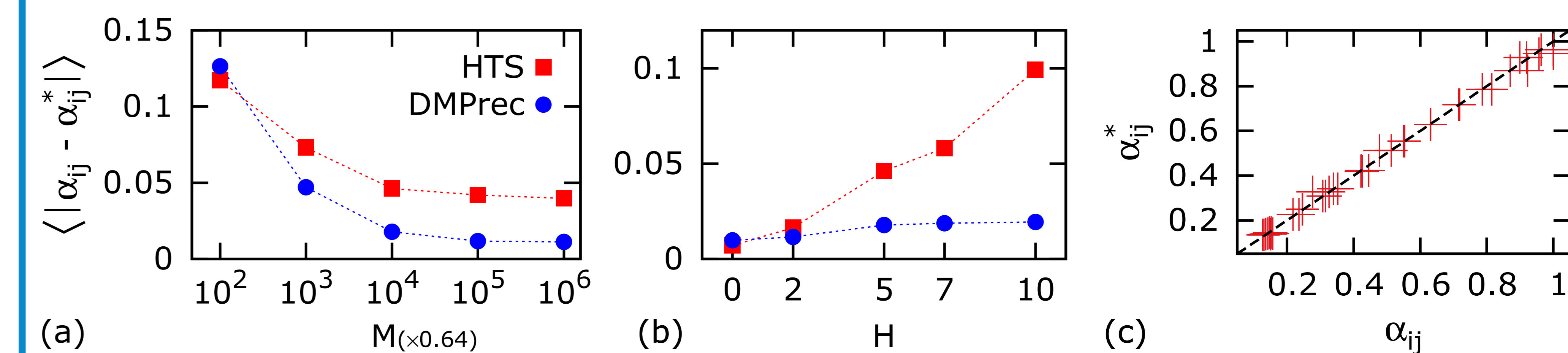
The DMPREC has a **much lower complexity** per step  $O(|E|^2 MT)$  vs.  $O(|E|^2 MNL_{H,T})$  for the HTS.

## CONSISTENCY OF DMPREC ON TREE NETWORKS

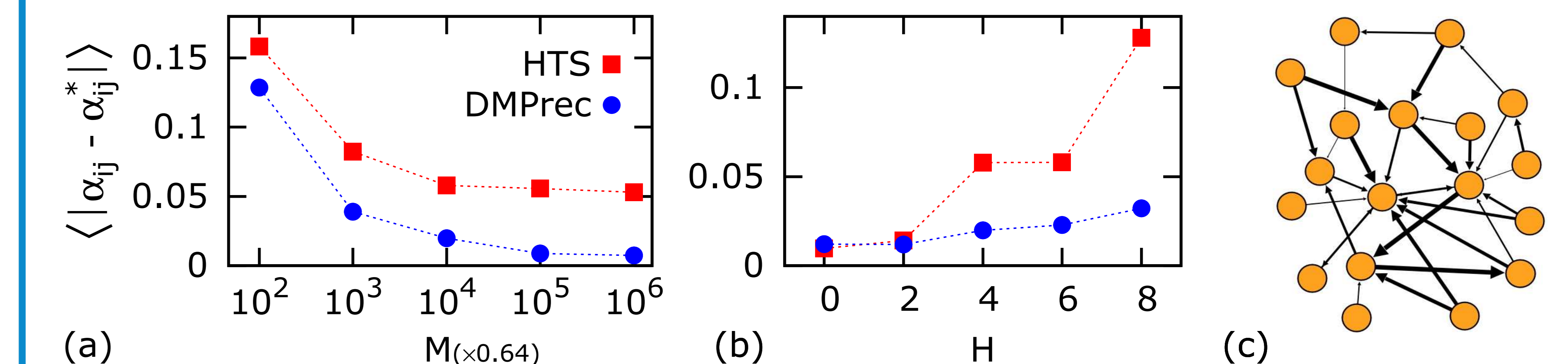
**Claim:** On tree graphs (regime in which DMP eqs. are derived),  $\lim_{M \rightarrow \infty} \frac{\partial f_{\text{DMP}}}{\partial \alpha_{rs}} |_{G_{\alpha^*}} = 0$ .  $\square$

## NUMERICAL RESULTS

**Comparison between HTS and DMPREC:**  
**days  $\rightarrow$  minutes**, for even **more accurate results**

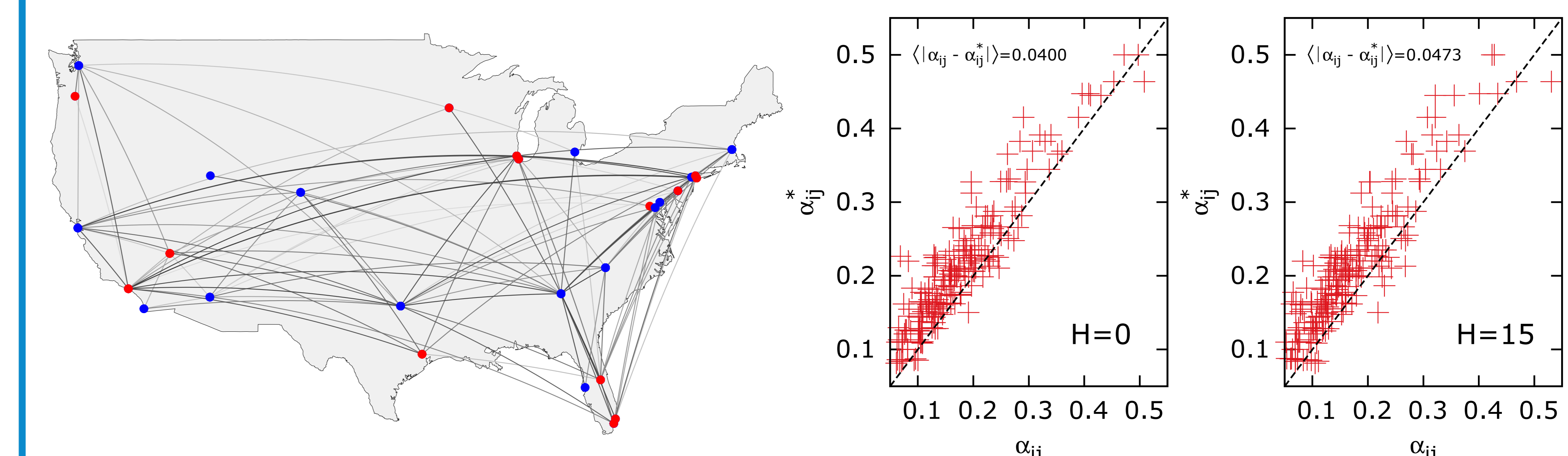


Small power-law network  $N = 20$  with random  $\alpha_{ij}^*$  in  $[0, 1]$  and  $T = 10$ : (a)  $H = 5$ , (b)  $M = 6400$ . Illustration for incomplete observations in time: (c)  $M = 6400$ , every other time stamp is missing.



Results for the network of relationships in a New England monastery: (a)  $H = 4$ , (b)  $M = 6400$ .

## Application to real data: air-traffic mediated epidemic spreading



Results for a sub-network of flights ( $|E| = 210$ ) between  $N = 30$  major U.S. hubs.  $M = 10,000$ .