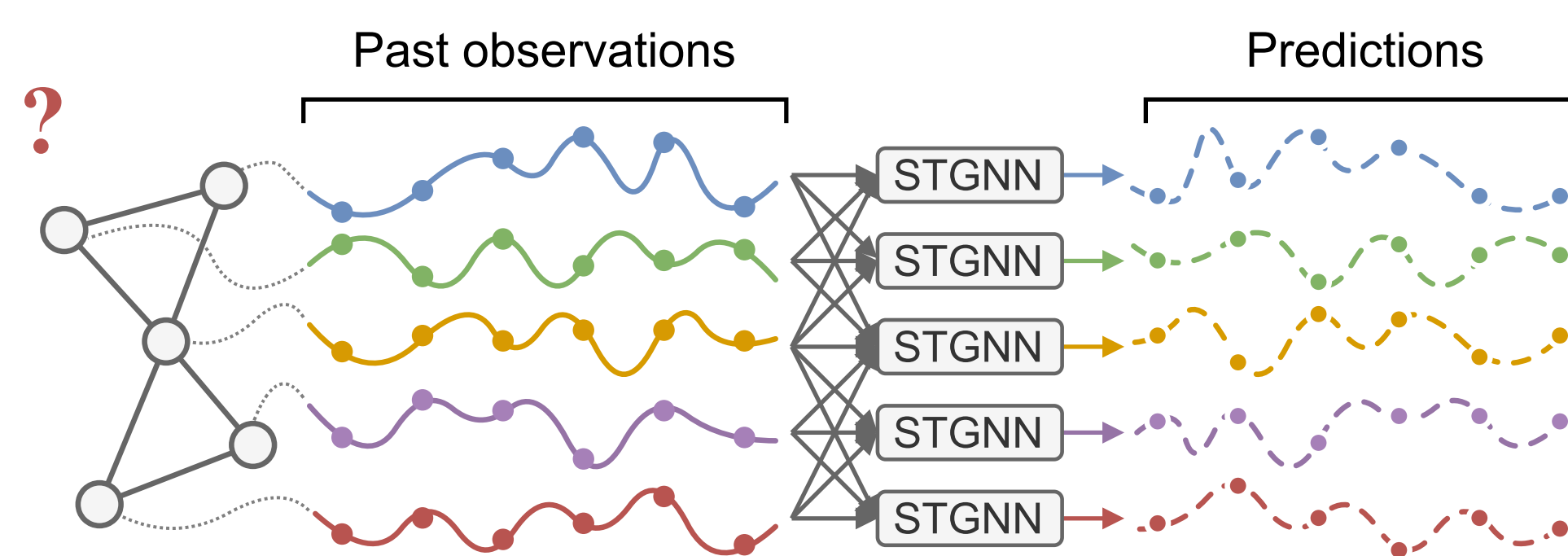


## Motivation

Spatiotemporal GNNs (STGNNs) are effective in forecasting time series collections, however, they rely on the available relational information.

- 😊 Relationships **localize predictions** w.r.t. related time series.
- 😞 Often, **no pre-defined graph** is accessible.
- 😞 Available **relational information** can be **incomplete** or misspecified.
- ❓ Can we **efficiently learn a graph** from data?
- ❓ Can we **keep computations sparse**?

## Latent Graph Learning



We want to **learn a graph** for spatiotemporal message-passing end-to-end.

- ! The learned graph should **maximize performance at task**.
- ! We want to **keep message-passing operations sparse**.
- ⚡ We exploit a **probabilistic framework** to learn **distributions over graphs**:

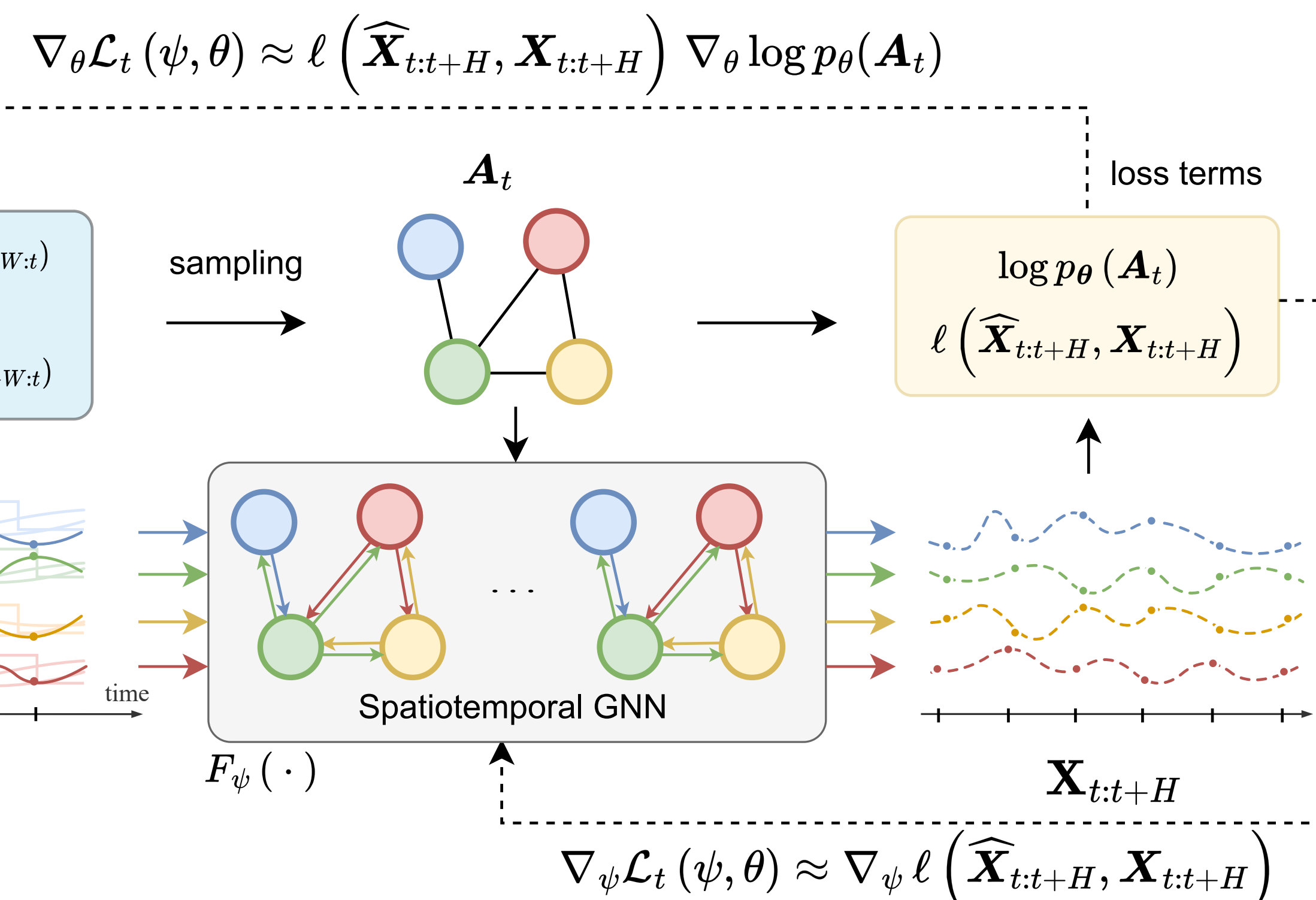
$$\widehat{\mathbf{A}}_t \sim p_\theta(\mathbf{A} | \mathcal{X}_{t-W:t}) \quad \widehat{\mathbf{X}}_{t:t+H} = F_\psi(\mathcal{X}_{t-W:t}, \widehat{\mathbf{A}}_t),$$

$$\hat{\theta}, \hat{\psi} = \arg \min_{\theta, \psi} \mathcal{L}_t(\psi, \theta) = \arg \min_{\theta, \psi} \mathbb{E}_{\widehat{\mathbf{A}}_t \sim p_\theta} \left[ \underbrace{\ell(\widehat{\mathbf{X}}_{t:t+H}, \mathbf{X}_{t:t+H})}_{\delta_t(\widehat{\mathbf{A}}_t; \psi)} \right]$$

## Score-based Gradient Estimators

$$\begin{aligned} \nabla_\theta \mathbb{E}_{p_\theta} [\delta_t(\widehat{\mathbf{A}}_t; \psi)] &= \mathbb{E}_{p_\theta} [\delta_t(\widehat{\mathbf{A}}_t; \psi) \nabla_\theta \log p_\theta(\widehat{\mathbf{A}}_t)] \\ &\approx \frac{1}{M} \sum_{i=1}^M \delta_t(\widehat{\mathbf{A}}_t^{(i)}; \psi) \nabla_\theta \log p_\theta(\widehat{\mathbf{A}}_t^{(i)}) \end{aligned}$$

- 😊 Computationally efficient.
- 😞 High variance gradient estimates.
- 😊 Allows for sparse computations.
- 😞 High sample complexity.



## Graph Samplers

## Binary Edge Sampler (BES)

$$\begin{aligned} p_\theta(\mathbf{A}_t[i, j] = 1) &= \\ &= \text{Bernoulli}(\sigma(\Phi_t[i, j])) \end{aligned}$$

- BES models the probability of sampling **each edge independently**.

## Subset Neighborhood Sampler (SNS)

$$\begin{aligned} p_\theta(\mathcal{N}(n) = S_K) &= \\ &= \sum_{\vec{S}_K \in \mathcal{P}(S_K)} \prod_{j \in \vec{S}_K} \frac{\exp(\Phi_t[n, j])}{1 - \sum_{k < j} \exp(\Phi_t[n, k])}. \end{aligned}$$

- SNS samples **K neighbors** for each node, imposing **sparsity** ( $|\mathcal{N}(n)| = K$ ).

## Improving Sample Efficiency

- ⚡ A **variance-reduced** estimator for each sampler (BES, SNS) based on control variates:

$$\nabla_\theta \mathbb{E}_{p_\theta} [\delta_t(\widehat{\mathbf{A}}_t; \psi)] \approx (\delta_t(\widehat{\mathbf{A}}_t; \psi) - \delta_t(\mathbf{A}_t^\mu; \psi)) \nabla_\theta \log p_\theta(\widehat{\mathbf{A}}_t).$$

We theoretically show that the **Fréchet mean** of the graph distribution is a sensible choice for  $\mathbf{A}_t^\mu$  (see Prep. 1–4).

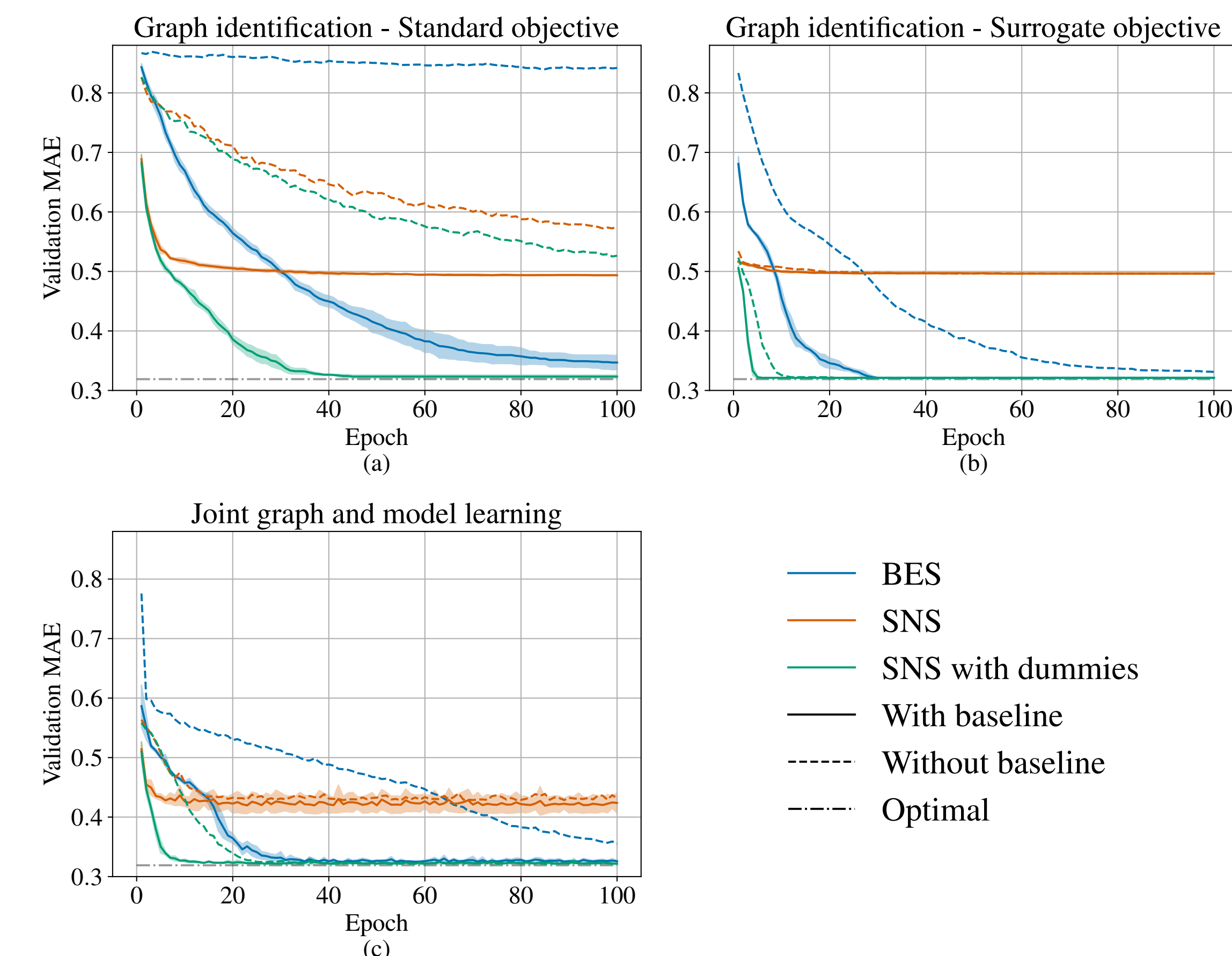
- ⚡ A **surrogate loss** to reweight the contribution of each error:

$$\nabla_\theta \mathcal{L}_t(\psi, \theta) \approx \mathbb{E}_{p_\theta} \left[ \lambda \delta_t(\widehat{\mathbf{A}}_t; \psi) \nabla_\theta \log p_\theta(\widehat{\mathbf{A}}_t) + \sum_{i=1}^N \delta_t^i(\widehat{\mathbf{A}}_t; \psi) \nabla_\theta \log p_\theta(\widehat{\mathbf{A}}_t[i, :]) \right],$$

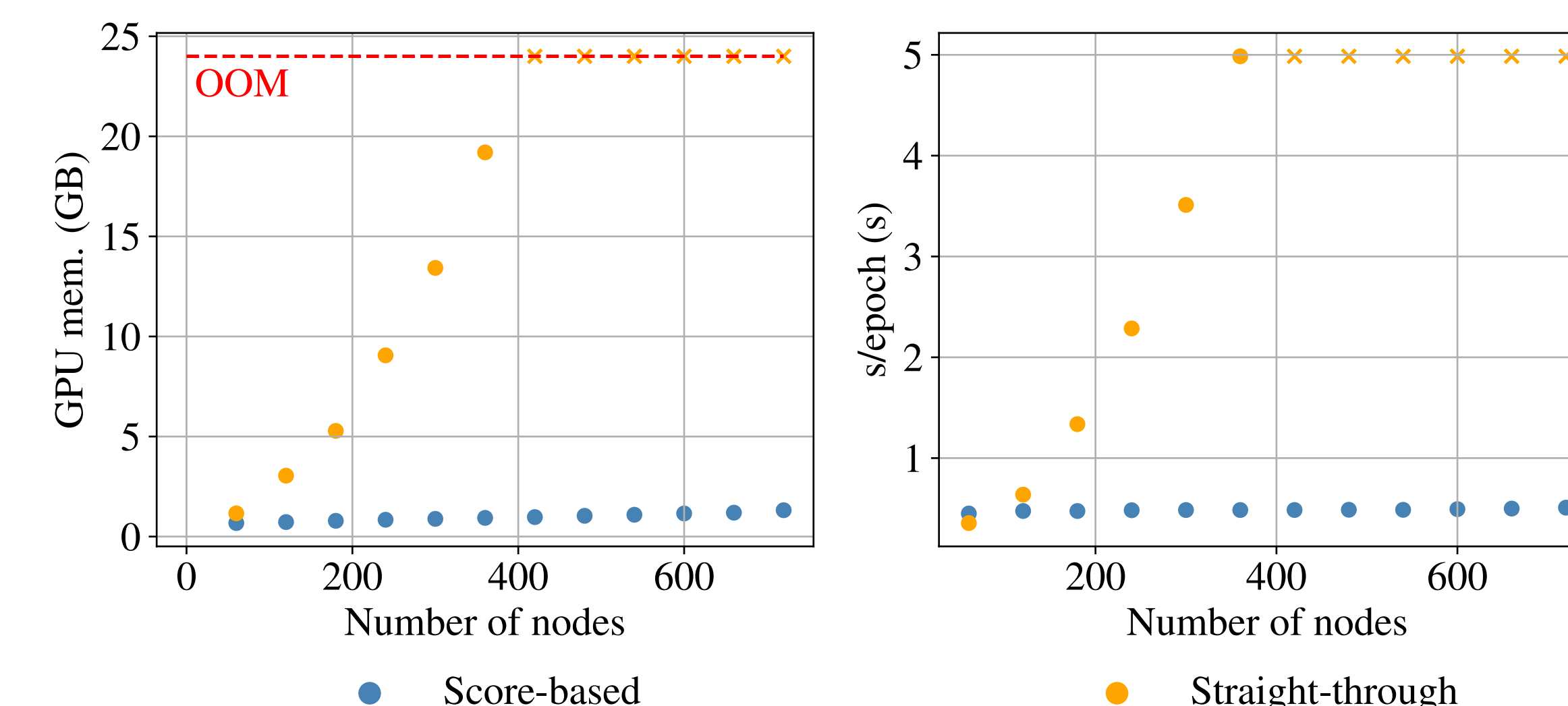
with  $\lambda$  trading of bias and variance (see Prep. 5 and Sec. 7.1). The surrogate objective **consistently reduces sample complexity**.

## Some Empirical Results

## Model training



## Computational scalability



\* Experiments on synthetic data.

andrea.cini@usi.ch

TorchSpatiotemporal/tsl