



Learning to Reconstruct Missing Data from Spatiotemporal Graphs with Sparse Observations

Ivan Marisca^{*1}

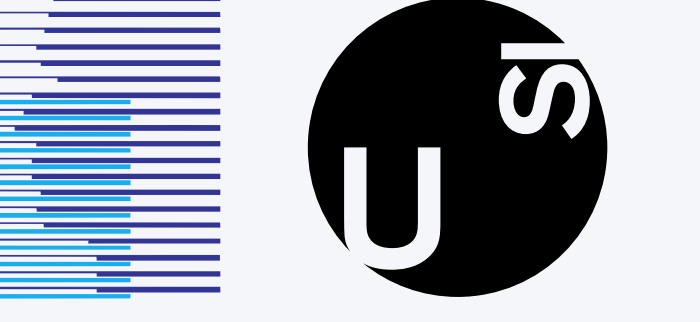
Andrea Cini^{*1}

Cesare Alippi¹²

¹The Swiss AI Lab IDSIA, Università della Svizzera italiana

²Politecnico di Milano

^{*}Equal contribution



Motivation

Missing values in spatiotemporal time series is an unavoidable problem when dealing with real-world applications.

- ▶ The state of the art includes **deep autoregressive methods**.
- ▶ Missing data blocks are often extended in time.
- ▶ Resulting observations are **highly sparse**.
- ▶ Autoregressive approaches suffer from **error compounding**.

Limits of Autoregressive Methods

SOTA autoregressive approaches model 3 different processes to account for the available data:

- ▶ $p(\mathbf{x}_t^i | \mathbf{X}_{<t})$, for **precedent** observations.
- ▶ $p(\mathbf{x}_t^i | \mathbf{X}_{>t})$, for **subsequent** observations.
- ▶ $p(\mathbf{x}_t^i | \{\mathbf{x}_t^{j \neq i}\})$, for **concurrent** observations.

DRAWBACKS:

- ▶ Computational **overhead**.
- ▶ **Error compounding**.
- ▶ Hard to capture **global context**.

- ▶ Drawbacks are more evident in **extremely sparse settings**.

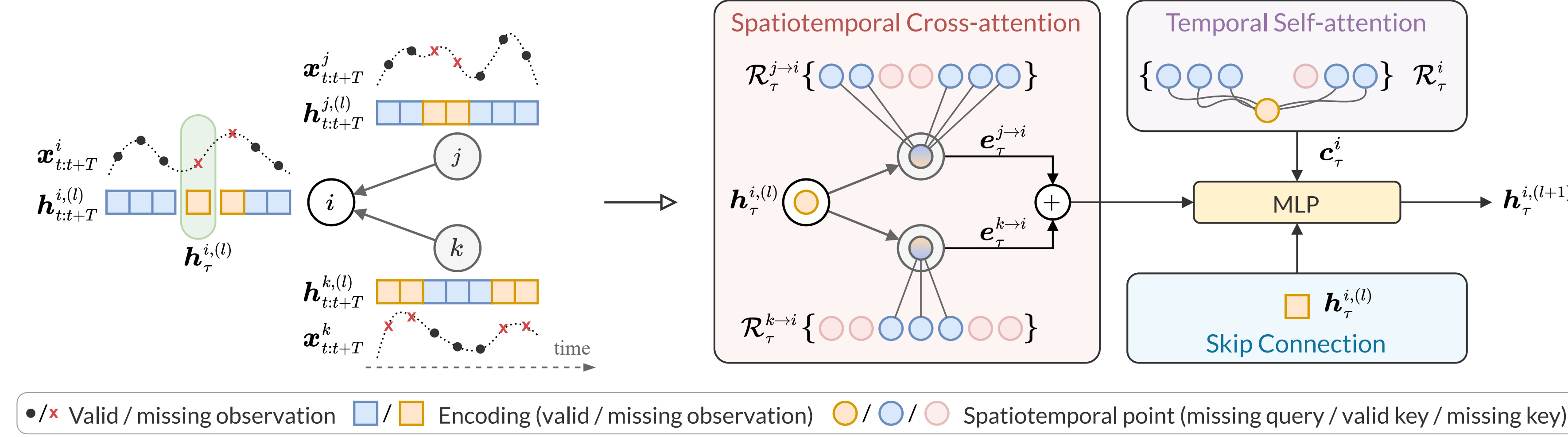
Our Different Approach

We associate **spatiotemporal coordinates** $\mathbf{Q}_t \in \mathbb{R}^{N \times d_q}$ to each point in time and space. And define:

- ▶ the **observed set** $\mathcal{X}_{t:t+T} = \{\langle \mathbf{x}_\tau^i, \mathbf{q}_\tau^i \rangle\}$ with valid observations.
- ▶ the **target set** $\mathcal{Y}_{t:t+T} = \{\mathbf{q}_\tau^i\}$ with target spatiotemporal points.

Given connectivity \mathbf{A} , for all $\mathbf{q}_\tau^i \in \mathcal{Y}_{t:t+T}$ we aim at modeling

$$p(\mathbf{x}_\tau^i | \mathbf{q}_\tau^i, \mathcal{X}_{t:t+T}, \mathbf{A}) \quad (1)$$



Spatiotemporal Point Inference Network (SPIN)

SPIN is an attention-based GNN imputation model f_θ such that

$$f_\theta(\mathbf{q}_\tau^i | \mathcal{X}_{t:t+T}, \mathbf{A}) \approx \mathbb{E}[p(\mathbf{x}_\tau^i | \mathbf{q}_\tau^i, \mathcal{X}_{t:t+T}, \mathbf{A})] \quad (2)$$

Representations at i -th node for each τ -th time step are learned with **two main components**:

TEMPORAL SELF-ATTENTION

- ▶ For each node \mathbf{h}_s^i associated with a **valid observation**, compute $\mathbf{r}_{s \rightarrow \tau}^i = \text{SelfMessage}(\mathbf{h}_s^i, \mathbf{h}_\tau^i)$
- ▶ Then compute **self-attention scores** α from \mathbf{r} and aggregate:

$$\mathbf{c}_\tau^i = \sum_s \alpha_{s \rightarrow \tau}^i \cdot \mathbf{r}_{s \rightarrow \tau}^i$$

SPATIOTEMPORAL CROSS-ATTENTION

- ▶ For each neighbor \mathbf{h}_s^j associated with a **valid observation**, we compute $\mathbf{r}_{s \rightarrow \tau}^{j \rightarrow i} = \text{CrossMessage}(\mathbf{h}_s^j, \mathbf{h}_\tau^i)$
- ▶ Then compute **cross-attention scores** α from \mathbf{r} and aggregate for each neighbor:

$$\mathbf{e}_\tau^{j \rightarrow i} = \sum_s \alpha_{s \rightarrow \tau}^{j \rightarrow i} \cdot \mathbf{r}_{s \rightarrow \tau}^{j \rightarrow i}$$

Then, target representation $\mathbf{h}_\tau^{i,(l)}$ is updated with a final aggregation step

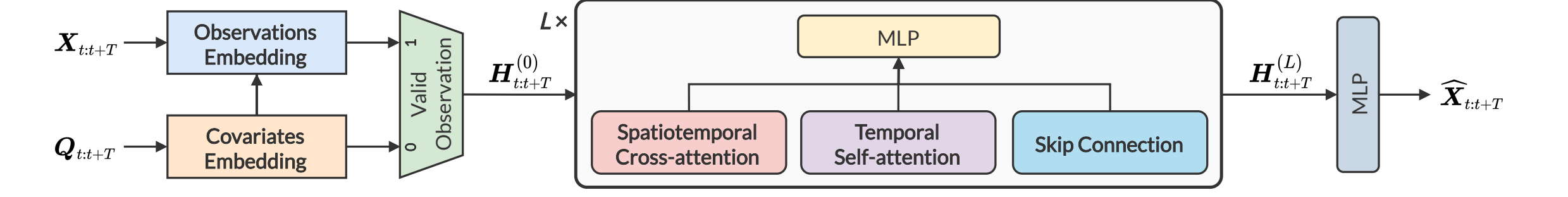
$$\mathbf{h}_\tau^{i,(l+1)} = \text{MLP}\left(\mathbf{h}_\tau^{i,(l)}, \mathbf{c}_\tau^i, \sum_{j \in \mathcal{N}(i)} \mathbf{e}_\tau^{j \rightarrow i,(l)}\right) \quad (3)$$

Imputations for all spatiotemporal points in $\mathcal{Y}_{t:t+T}$ are obtained as

$$\hat{\mathcal{Y}}_{t:t+T} = \{\hat{\mathbf{x}}_\tau^i = \text{MLP}(\mathbf{h}_\tau^{i,(L)}) | \mathbf{q}_\tau^i \in \mathcal{Y}_{t:t+T}\} \quad (4)$$

Covariates are **positional encodings** for observations: representation $\mathbf{h}_\tau^{i,(0)}$ is initialized as

$$\mathbf{h}_\tau^{i,(0)} = \text{MLP}(\mathbf{q}_\tau^i, \mathbf{x}_\tau^i) \quad \text{if } \langle \mathbf{x}_\tau^i, \mathbf{q}_\tau^i \rangle \in \mathcal{X}_{t:t+T} \quad \mathbf{h}_\tau^{i,(0)} = \text{MLP}(\mathbf{q}_\tau^i) \quad \text{if } \mathbf{q}_\tau^i \in \mathcal{Y}_{t:t+T}$$



The SPIN architecture.

Hierarchical Attention

The base SPIN layer has $\mathcal{O}((N + E)T^2)$ complexity. To remove the quadratic term, we rewire attention to be **hierarchical**:

- ▶ Add K **dummy nodes** that act as hubs.
- ▶ Update the $\mathbf{z}^i \in \mathbb{R}^{d_z}$ hubs' representations by **querying the available nodes** $\{\mathbf{h}_\tau^i\}$.
- ▶ Update node encoding \mathbf{h}_τ^i by **querying the updated hubs** $\{\tilde{\mathbf{z}}_k^j\}$ for each neighbor $j \in \mathcal{N}(i)$.

This reduces complexity to $\mathcal{O}((N + E)KT)$ with $K \ll T$. We refer to this variation as **SPIN-H**.

Some Empirical Results

We report here the performance in reconstructing long blocks of consecutive missing values (from 12 to 36 time steps).

Table 1. Performance (MAE) with an increasing number of simulated failures.

	METR-LA			PEMS-BAY		
	Failure probability 5 % 10 % 15 %			Failure probability 5 % 10 % 15 %		
BRITS	5.87 ± 0.03	7.26 ± 0.06	8.29 ± 0.07	4.14 ± 0.05	5.41 ± 0.08	5.84 ± 0.04
SAITS	4.73 ± 0.07	6.66 ± 0.05	7.27 ± 0.03	3.88 ± 0.09	7.62 ± 0.21	8.01 ± 0.11
Transformer	6.03 ± 0.04	7.19 ± 0.05	8.06 ± 0.05	3.69 ± 0.06	5.09 ± 0.05	6.02 ± 0.04
GRIN	3.05 ± 0.02	4.52 ± 0.05	5.82 ± 0.06	2.26 ± 0.03	3.45 ± 0.06	4.35 ± 0.04
SPIN	2.71 ± 0.02	3.32 ± 0.02	3.87 ± 0.05	1.78 ± 0.03	2.15 ± 0.03	2.41 ± 0.02
SPIN-H	2.64 ± 0.02	3.17 ± 0.02	3.61 ± 0.04	1.75 ± 0.04	2.16 ± 0.03	2.48 ± 0.02

Our library for neural spatiotemporal data processing:

TorchSpatiotemporal/tsl