

# Scalable Spatiotemporal Graph Neural Networks

Andrea Cini<sup>\*1</sup>

Ivan Marisca<sup>\*1</sup>

Filippo Maria Bianchi<sup>23</sup>

Cesare Alippi<sup>14</sup>

<sup>1</sup>The Swiss AI Lab IDSIA, Università della Svizzera italiana

<sup>2</sup>UiT the Arctic University of Norway

<sup>3</sup>NORCE Norwegian Research Centre

<sup>4</sup>Politecnico di Milano

<sup>\*</sup>Equal contribution

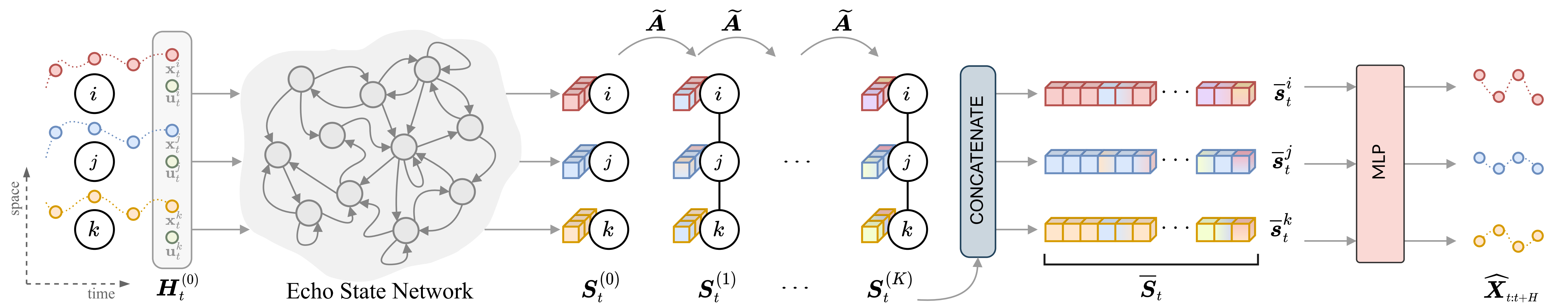
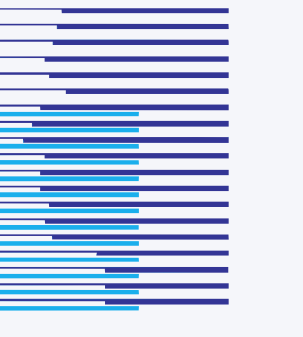


Figure 1. Scalable Graph Predictor

## SGP: Scalable Graph Predictor

A novel approach based on an **encode-decoder** architecture with a **training-free spatiotemporal encoding** scheme and where the only learned parameters are in the **node-level trainable decoder** (an MLP).

- ▶ Representations for each point in time and space can be **precomputed**.
- ▶ The decoder can be trained by **sampling uniformly time and space**.
- ▶ This enables **scalability at training time**:  $\mathcal{O}(ET) \rightarrow \mathcal{O}(1)$ .

## Motivation

**Scalability in Spatiotemporal GNNs** is challenging to achieve due to the large amounts of input data.

- ▶ Standard STGNNs have **time and space complexity** of  $\mathcal{O}(ET)$ .
- ▶ If an attention component is included, **complexity becomes quadratic** in either time or space.
- ▶ **Subsampling strategies** can break long-range spatiotemporal dependencies and are prone to failure.

## Spatiotemporal Encoder

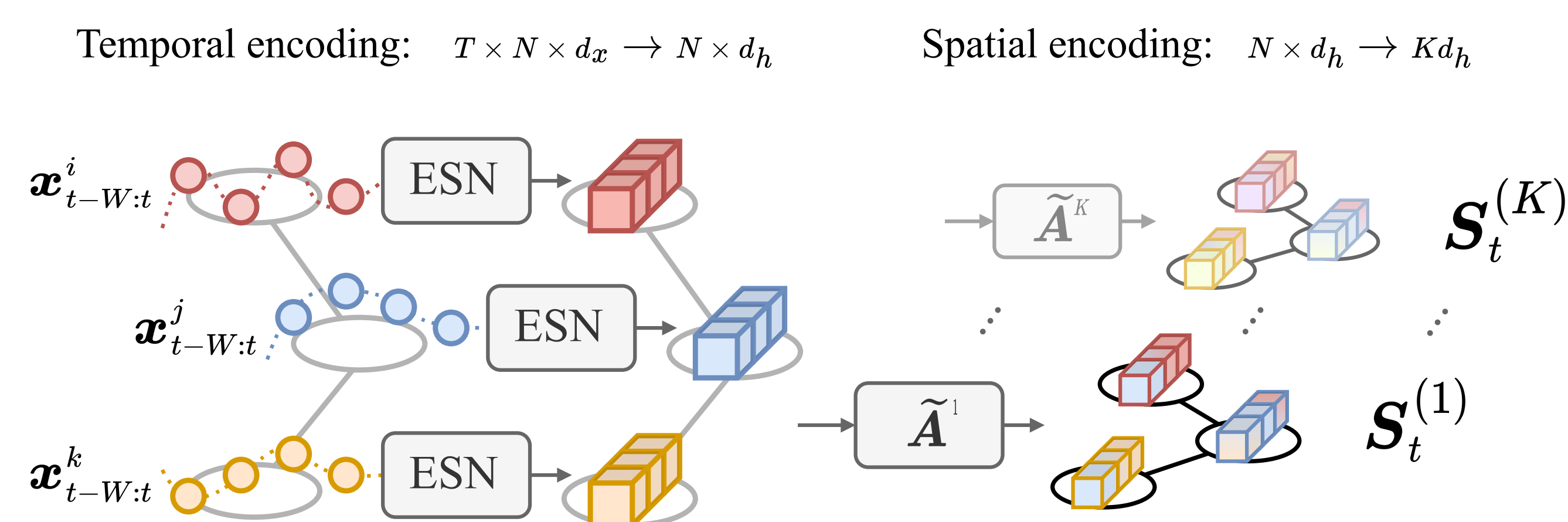


Figure 2. Spatiotemporal encoding scheme.

The spatiotemporal encoder relies on **two modules**.

- ▶ A **randomized recurrent neural network** for encoding sequences.
- ▶ A **propagation process** through the graph structure.

### TEMPORAL ENCODING

We adopt the **deep echo state network** framework [1].

$$\begin{aligned} \mathbf{h}_t^{i,(0)} &= [\mathbf{x}_t^i \| \mathbf{u}_t^i], \\ \hat{\mathbf{h}}_t^{i,(l)} &= \tanh(\mathbf{W}_u^{(l)} \mathbf{h}_t^{i,(l-1)} + \mathbf{W}_h^{(l)} \mathbf{h}_{t-1}^{i,(l)} + \mathbf{b}^{(l)}), \\ \mathbf{h}_t^{i,(l)} &= (1 - \gamma_l) \mathbf{h}_{t-1}^{i,(l)} + \gamma_l \hat{\mathbf{h}}_t^{i,(l)}. \quad l = 1, \dots, L \end{aligned}$$

- ▶ Weight matrices are **randomly generated**.
- ▶ Proper normalization makes the system **stable**.
- ▶ The concatenated multilayer state representations **encode rich, multi-scale, dynamics**

$$\bar{\mathbf{H}}_t = (\mathbf{H}_t^{(0)} \| \mathbf{H}_t^{(1)} \| \dots \| \mathbf{H}_t^{(L)}).$$

### SPATIAL PROPAGATION

The extracted temporal encodings are propagated by using **powers of a graph shift operator**.

$$\begin{aligned} \mathbf{S}_t^{(0)} &= \bar{\mathbf{H}}_t = (\mathbf{H}_t^{(0)} \| \mathbf{H}_t^{(1)} \| \dots \| \mathbf{H}_t^{(L)}), \\ \mathbf{S}_t^{(k)} &= \tilde{\mathbf{A}}^k \mathbf{S}_t^{(k-1)} = (\tilde{\mathbf{A}}^k \mathbf{H}_t^{(0)} \| \tilde{\mathbf{A}}^k \mathbf{H}_t^{(1)} \| \dots \| \tilde{\mathbf{A}}^k \mathbf{H}_t^{(L)}), \\ \bar{\mathbf{S}}_t &= (\mathbf{S}_t^{(0)} \| \mathbf{S}_t^{(1)} \| \dots \| \mathbf{S}_t^{(K)}). \end{aligned}$$

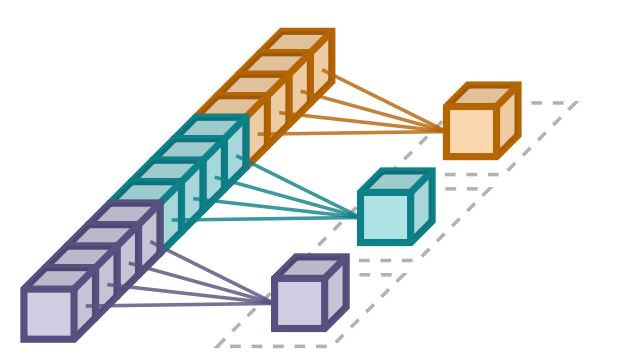
- ▶ This can be done either **recursively** or **in parallel**.

## Multi-Scale Decoder

Could be a standard MLP, however, we can exploit **the structure of the embedding**.

- ▶ We make the connectivity of the first MLP layer **sparse**.

$$\begin{aligned} \mathbf{Z}_t^{(k)} &= \sigma(\tilde{\mathbf{A}}^k \mathbf{H}_t^{(0)} \Theta_k^{(0)} \| \dots \| \tilde{\mathbf{A}}^k \mathbf{H}_t^{(L)} \Theta_k^{(L)}) \\ &= \sigma\left(\mathbf{S}_t^{(k)} \begin{bmatrix} \Theta_k^{(0)} & \dots & 0 \\ 0 & \dots & \Theta_k^{(L)} \end{bmatrix}\right), \\ \bar{\mathbf{Z}}_t &= (\mathbf{Z}_t^{(0)} \| \mathbf{Z}_t^{(1)} \| \dots \| \mathbf{Z}_t^{(K)}). \end{aligned}$$



- ▶ Standard fully connected layers can then be used to obtain predictions.
- ▶ Trained by sampling **minibatches of representations** over time and space **as if they were iid**.

## Some Results

We tested the scalability of the approach by imposing a limit on training time (1 hour) and capping GPU memory utilization (12 GB).

		Prediction error (MAE)			Resource utilization		
		30 mins	7 hours 30 mins	11 hours	Batch/s	Memory	Batch size
PV-US 100	DCRNN	1.39 ± 0.09	3.34 ± 0.22	<b>3.54 ± 0.48</b>	2.04 ± 0.01	9.63 GB	2
	GWNet	1.45 ± 0.13	5.09 ± 0.63	5.26 ± 1.34	2.01 ± 0.02	11.64 GB	2
	GatedGN	1.33 ± 0.08	<b>2.94 ± 0.05</b>	<b>3.12 ± 0.14</b>	8.41 ± 0.09	11.46 GB	5
PV-US Full	<b>SGP</b>	<b>1.09 ± 0.01</b>	<b>3.14 ± 0.21</b>	<b>3.16 ± 0.19</b>	<b>116.58 ± 8.74</b>	<b>2.21 GB</b>	4096
PV-US Full	DCRNN	1.59 ± 0.17	4.10 ± 0.27	4.93 ± 0.60	1.37 ± 0.00	11.59 GB	1*
	GWNet	1.65 ± 0.23	6.93 ± 0.58	7.93 ± 0.17	0.77 ± 0.00	11.35 GB	2
	GatedGN	1.61 ± 0.06	3.25 ± 0.04	<b>3.04 ± 0.05</b>	8.83 ± 0.10	11.14 GB	1*
PV-US Full	<b>SGP</b>	<b>1.09 ± 0.00</b>	<b>3.06 ± 0.11</b>	<b>3.13 ± 0.13</b>	<b>118.64 ± 8.35</b>	<b>2.21 GB</b>	4096

Table 1. Results on large-scale datasets.

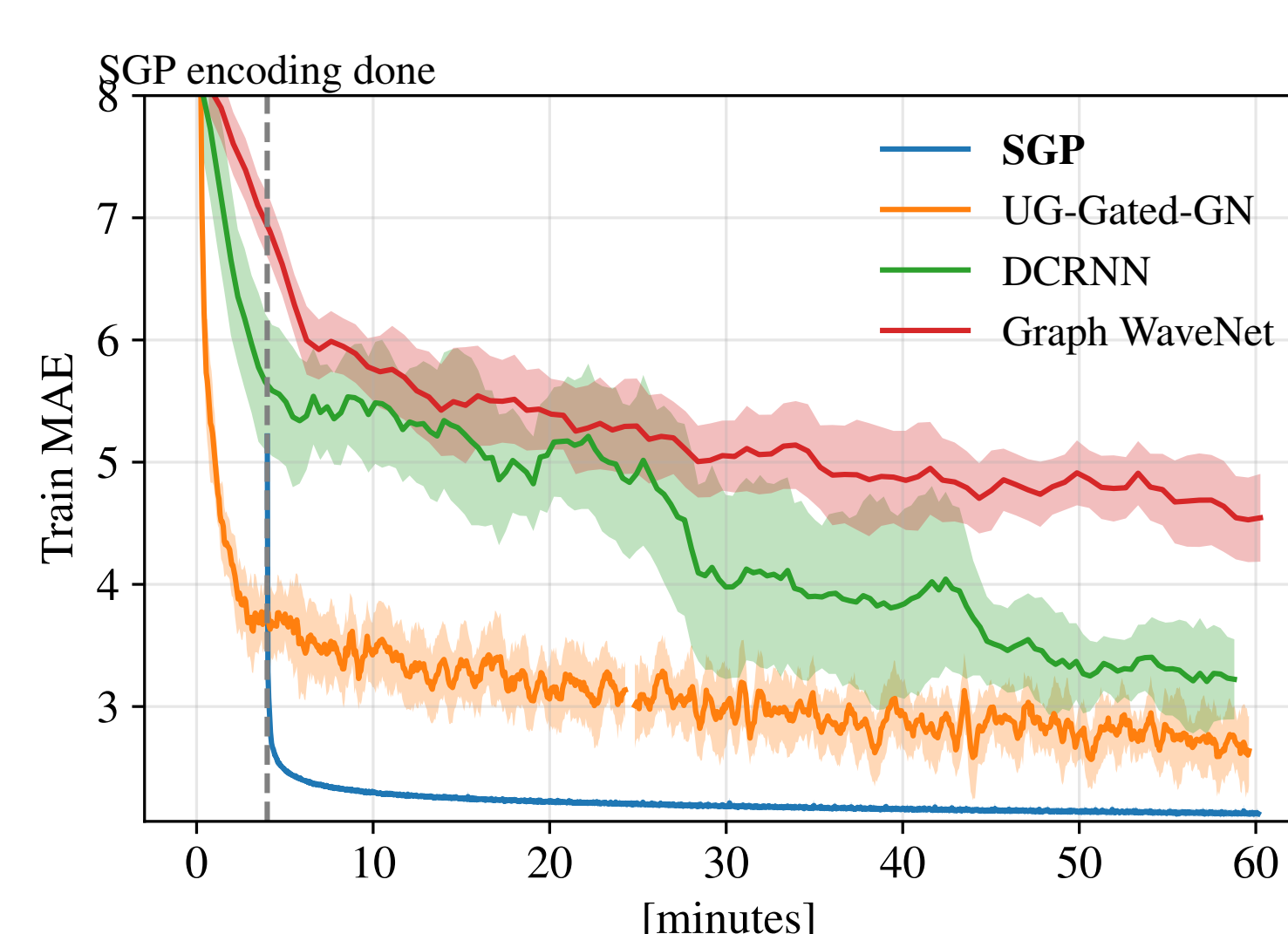


Figure 3. Learning curve on PV-US.

## References

- [1] C. Gallicchio, A. Micheli, and L. Pedrelli. Deep reservoir computing: A critical experimental analysis. *Neurocomputing*, 268: 87–99, 2017.

Our library for spatiotemporal data processing:

**TorchSpatiotemporal/tsl**