

Learning on Graphs Conference · November 28, 2024

# Graph Deep Learning for Time Series Processing

Forecasting, Reconstruction and Analysis

---

Andrea **Cini**, Ivan **Marisca**, Daniele **Zambon**

Graph Machine Learning Group ([gmlg.ch](https://gmlg.ch))

The Swiss AI Lab IDSIA

Università della Svizzera italiana



idsia



# Introduction

---



Traffic monitoring



Smart cities



Energy analytics



Physics

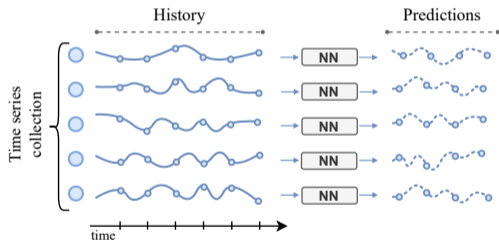


Stock markets

# Deep learning for time series forecasting

Modern deep learning forecasting methods rely on a **single neural network** trained on a collection of **related time series**.

- 😊 Each time series is processed **independently**.
- 😊 Parameters are **shared**.
- 😊 Effective and **sample efficient**.
- 😞 **Dependencies are neglected**.



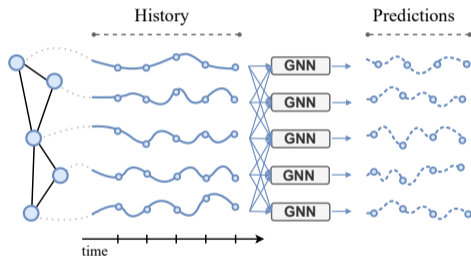
[1] Salinas *et al.*, "DeepAR: Probabilistic forecasting with autoregressive recurrent networks", IJF 2020.

[2] Benidis *et al.*, "Deep Learning for Time Series Forecasting: Tutorial and Literature Survey", ACM CS 2022.

# Graph deep learning for time series forecasting

We will show **graph deep learning (GDL)** provides appropriate operators to **go beyond these limitations**.

- 😊 **Dependencies** are **embedded into the processing** as inductive biases.
- 😊 Operate on **sets of correlated time series**.
- 😊 Parameters are **shared**.



☹️ There are inherent **challenges** in applying this processing to data from the real world.

# What this tutorial is about

---

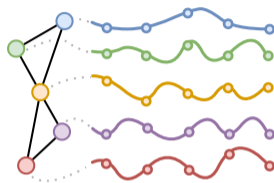
This tutorial presents advances coming from the **combination** of

1. **deep learning** for **time series** and
2. **deep learning** on **graphs**.

The **objective** of the tutorial is to provide:

1. a comprehensive framework for **graph-based time series processing** models;
2. methods to address **challenges** and potential **pitfalls**;
3. **tools** and **guidelines** for **real-world applications** and developing **new methods**.

This presentation is complemented by a **demo** and a **tutorial paper** [3].

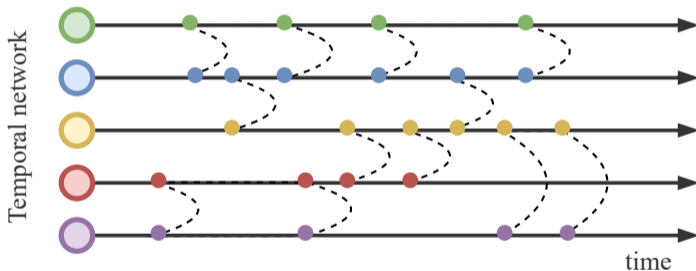


---

[3] Cini, Marisca, Zambon, and Alippi, “Graph Deep Learning for Time Series Forecasting”, Preprint 2023.

# What this tutorial is not about

⚠ This tutorial is **not** about processing sequences of interactions in **temporal networks**.



→ Graphs will be a representation of the (dynamic) **relationships among** (possibly irregular) **time series**.


# Tutorial outline

---

## Part 1

- 1.1)** Correlated time series
- 1.2)** Graph-based representation
- 1.3)** STGNN architectures
- 1.4)** Global and local models
-  Software demo

## Part 2

- 2.1)** Scalability
- 2.2)** Dealing with missing data
- 2.3)** Latent graph learning
- 2.4)** Model quality assessment
-  Conclusions



Part 1

# Graph-based Processing of Correlated Time series

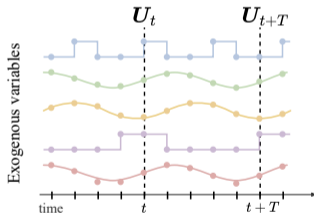
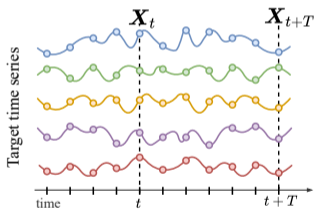
## **Correlated time series**

---

# Collections of time series

We consider a set  $\mathcal{D}$  of  $N$  **correlated time series**. Each  $i$ -th time series can be associated with:

- **observations**  $x_t^i \in \mathbb{R}^{d_x}$  at each time step  $t$ ;
- **exogenous variables**  $u_t^i \in \mathbb{R}^{d_u}$  at each time step  $t$ ;
- a vector of **static (time-independent) attributes**  $v^i \in \mathbb{R}^{d_v}$ .



Static attributes



Capital letters denote the stacked  $N$  time series, i.e.,  $X_t \in \mathbb{R}^{N \times d_x}$ ,  $U_t \in \mathbb{R}^{N \times d_u}$ .

→ We call **spatial** the dimension spanning the collection.

# Correlated time series

We consider a **time-invariant** stochastic process generating each time series as

$$\mathbf{x}_t^i \sim p^i \left( \mathbf{x}_t^i | \mathbf{X}_{<t}, \mathbf{U}_{\leq t}, \mathbf{V} \right) \quad \text{for all } i = 1 \dots N, t = 0, \dots, T - 1$$

and assume the existence of a **causality à la Granger** among time series.

Furthermore time series

- are assumed
  - a) homogenous, b) synchronous, c) regularly sampled.
- can be generated by different processes.

**Notation:**

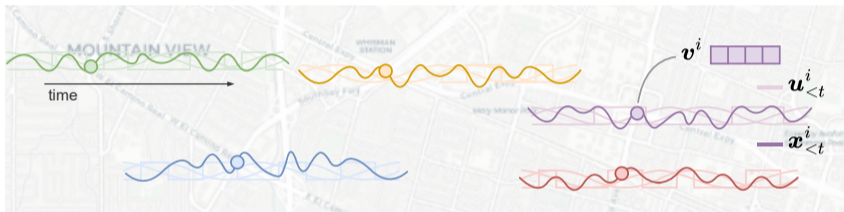
$$\mathcal{X}_t = \langle \mathbf{X}_t, \mathbf{U}_t, \mathbf{V} \rangle$$

$$\mathcal{X}_{<t} = [\mathcal{X}_0, \dots, \mathcal{X}_{t-2}, \mathcal{X}_{t-1}]$$

! Assumptions a),b),c) can be relaxed as we will discuss in the 2nd part.

## Example: Traffic monitoring system

Consider a sensor network monitoring the speed of vehicles at crossroads.



- $\mathbf{X}_{<t}$  collects past traffic speed measurements.
- $\mathbf{U}_t$  stores identifiers for time-of-the-day and day-of-the-week.
- $\mathbf{V}$  collects static sensor's features, e.g., type or number of lanes of the monitored road.

→ Strong dependencies among time series that reflect the road network.

# Forecasting

---

# Time series forecasting

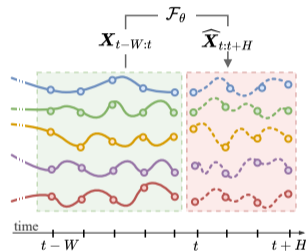
## Multi-step time-series forecasting

Given a window of  $W \geq 1$  past values

$$\mathcal{X}_{t-W:t} = [\mathcal{X}_{t-W}, \dots, \mathcal{X}_{t-1}],$$

predict  $H \geq 1$  future observations

$$\mathbf{X}_{t+h} \quad h = 1, \dots, H.$$



In particular, we are interested in learning a parametric model  $\mathcal{F}(\cdot; \theta)$  s.t.

$$\mathcal{F}(\mathcal{X}_{t-W:t}, \mathbf{U}_{t:t+H}; \theta) = \widehat{\mathbf{X}}_{t:t+H} \approx E_p[\mathbf{X}_{t:t+H}].$$

Probabilistic predictors can be considered as well, but we focus on point forecasts.

# Training objective

---

For point predictors, parameters  $\theta$  can be learned by **minimizing a cost function**  $\ell(\cdot, \cdot)$  (e.g., MSE) on a training set

$$\begin{aligned}\hat{\theta} &= \arg \min_{\theta} \frac{1}{NT} \sum_{t=1}^T \ell \left( \widehat{\mathbf{X}}_{t:t+H}, \mathbf{X}_{t:t+H} \right) \\ &= \arg \min_{\theta} \frac{1}{NT} \sum_{t=1}^T \left\| \mathbf{X}_{t:t+H} - \widehat{\mathbf{X}}_{t:t+H} \right\|_2^2.\end{aligned}$$

- ! Choosing a different cost function allows for predicting other values.  
→ **Example:** minimizing the MAE results in forecasts of the median.



# Global and local predictors

## Local models



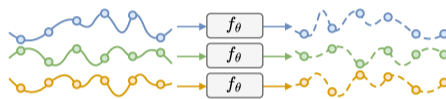
$$\hat{x}_{t+h}^i = f(x_{t-W:t}^i; \theta^i)$$

**Example:** Box-Jenkins method

😊 Tailored to each time series.

☹ Inefficient.

## Global models



$$\hat{x}_{t+h}^i = f(x_{t-W:t}^i; \theta)$$

**Example:** DeepAR [1]

😊 Sample efficient.

😊 Allows for more complex models.

☹ Both approaches neglect dependencies among time series.

[1] Salinas *et al.*, “DeepAR: Probabilistic forecasting with autoregressive recurrent networks”, IJF 2020.

[4] Montero-Manso *et al.*, “Principles and algorithms for forecasting groups of time series: Locality and globality”, IJF 2021.

# Accounting for spatial dependencies

- One option is to consider the input as **single multivariate time series**

→ Resulting predictors are **local**:  $\widehat{\mathbf{X}}_{t+h} = f(\mathbf{X}_{t-W:t}, \dots; \boldsymbol{\theta})$ .

☹ High **sample complexity** and poor **scalability**.

- Models **operating on sets of time series** would allow to keep parameters shared.

→ Resulting predictors are **global**:  $\widehat{\mathbf{X}}_{t+h}^S = \mathcal{F}(\mathbf{X}_{t-W:t}^S, \dots; \boldsymbol{\theta}), \quad \forall S \subseteq \mathcal{D}$

😊 Can be implemented by **attention-based** models (e.g, **Transformers**).

☹ **Does not exploit structural priors**, **high computational** and **sample complexity**.

- Other methods (e.g., [5]) rely on **dimensionality reduction** to extract **shared latent factors**.

😊 Might work well if data are **low-rank**.

☹ **Local** and **relational information** are **lost** and can still suffer from, **scalability** issues.

[2] Benidis *et al.*, “Deep Learning for Time Series Forecasting: Tutorial and Literature Survey”, ACM CS 2022.

[5] Sen *et al.*, “Think globally, act locally: A deep neural network approach to high-dimensional time series forecasting”, NeurIPS 2019.

# **Graph-based representation**

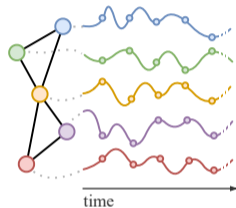
---

# Relational information

💡 Exploit **functional dependencies** as an **inductive bias** to improve the forecasts.

We can model pairwise relationships existing at time step  $t$  with **adjacency matrix**  $\mathbf{A}_t \in \{0, 1\}^{N \times N}$ .

- $\mathbf{A}_t$  can be **asymmetric** and **dynamic** (can vary with  $t$ ).

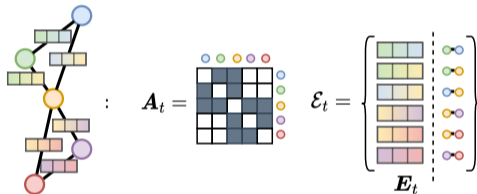


# Relational information with attributes

Optional **edge attributes**  $e_t^{ij} \in \mathbb{R}^{d_e}$  can be associated to each non-zero entry of  $\mathbf{A}_t$ .

The **set of attributed edges** is denoted by

$$\mathcal{E}_t \doteq \{ \langle (i, j), e_t^{ij} \rangle \mid \forall i, j : \mathbf{A}_t[i, j] \neq 0 \}.$$



→ Edge attributes can be both **categorical** or **numerical**.

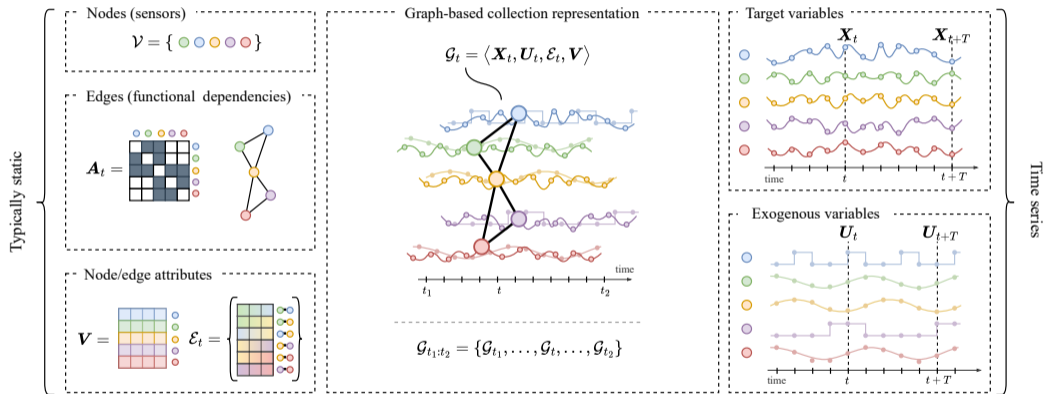
# Example: Traffic monitoring system

Consider again the sensor network of the previous example.



- Edges in  $\mathcal{E}$  can be obtained by considering the **road network**.
  - Road closures and traffic diversions can be accounted for with a dynamic topology  $\mathcal{E}_t$ .

# Graph-based representations for correlated time series



$\mathcal{G}_t \doteq \langle \mathbf{X}_t, \mathbf{U}_t, \mathcal{E}_t, \mathbf{V} \rangle$  contains the available information w.r.t. time step  $t$ .

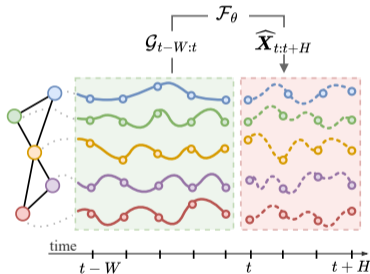
# Relational inductive biases for time series forecasting

Forecasts can be conditioned on the available relational information  $\mathcal{E}_{t-W:t}$

$$\widehat{\mathbf{X}}_{t:T+H}^S = \mathcal{F} \left( \mathcal{G}_{t-W:t}^S, \mathbf{U}_{t:t+H}^S; \boldsymbol{\theta} \right) \quad \forall S \in \mathcal{D}$$

The conditioning can act as a **regularization** to localize predictions w.r.t. **each node**.

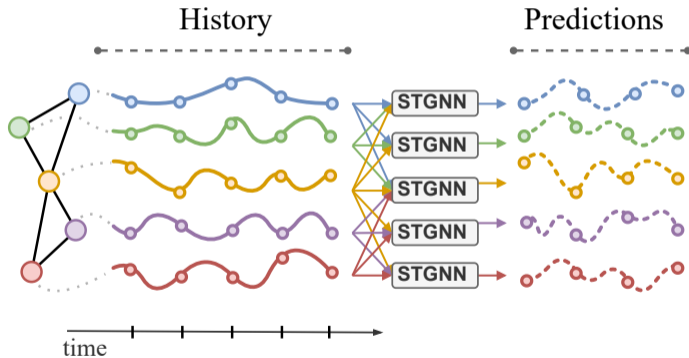
- 😊 Relational priors **prune spurious correlations**.
- 😊 More **scalable** than standard multivariate models.
- 😊 Can forecast any **subset** of correlated time series.





# Spatiotemporal graph neural networks

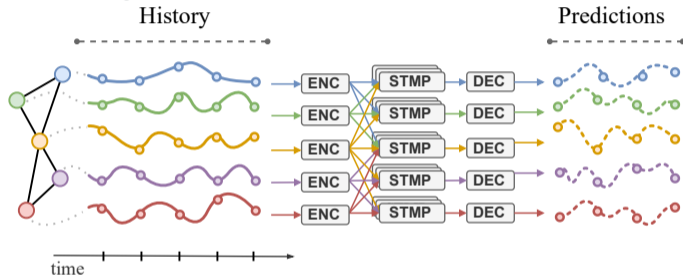
We call [spatiotemporal graph neural networks \(STGNNs\)](#) a neural network exploiting both temporal and spatial relations of the input spatiotemporal time series.



We focus on models based on [message passing \(MP\)](#).

# A general recipe for building STGNNs

We consider STGNNs consisting of three main components



- $ENC(\cdot)$  is the **encoding** layer, e.g., implemented by an MLP.
- $STMP(\cdot)$  is a stack of **spatiotemporal message-passing (STMP)** layers.
- $DEC(\cdot)$  is the **readout** layer, e.g., implemented by an MLP.

## A closer look

---

Representations are updated as follows.

$$\mathbf{h}_{t-1}^{i,0} = \text{ENCODER} \left( \mathbf{x}_{t-1}^i, \mathbf{u}_{t-1}^i, \mathbf{v}^i \right), \quad (1)$$

$$\mathbf{H}_{t-1}^{l+1} = \text{STMP}^l \left( \mathbf{H}_{\leq t-1}^l, \mathcal{E}_{\leq t-1} \right), \quad l = 0, \dots, L-1 \quad (2)$$

$$\hat{\mathbf{x}}_{t:t+H}^i = \text{DECODER} \left( \mathbf{h}_{t-1}^{i,L}, \mathbf{u}_{t:t+H}^i \right). \quad (3)$$

- **ENC**( · ) process each observation **independently**.
- **STMP**( · ) is where **propagation** through time and space happens.
- **DEC**( · ) maps each representation to **predictions**.

---

[3] Cini *et al.*, “Graph Deep Learning for Time Series Forecasting”, Preprint 2023.

# Spatiotemporal message-passing (STMP)

STMP blocks can be defined as:

$$\mathbf{h}_t^{i,l+1} = \text{UP}^l \left( \mathbf{h}_{\leq t}^{i,l}, \text{AGGR}_{j \in \mathcal{N}_t(i)} \left\{ \text{MSG}^l(\mathbf{h}_{\leq t}^{i,l}, \mathbf{h}_{\leq t}^{j,l}, \mathbf{e}_{\leq t}^{ji}) \right\} \right)$$

Each block processes **sequences** while accounting for **relational dependencies**.

As in standard MP operators:

- $\text{MSG}^l(\cdot)$  is a **message function**, e.g., implemented by *temporal convolutional layers*.
- $\text{AGGR}\{\cdot\}$  is a permutation invariant **aggregation function**.
- $\text{UP}^l(\cdot)$  is an **update function**, e.g., implemented by an RNN.

! Blocks can be implemented by composing MP and sequence modeling operators.

→ Many possible designs exist.

---

[3] Cini *et al.*, “Graph Deep Learning for Time Series Forecasting”, Preprint 2023.

[6] Gilmer *et al.*, “Neural message passing for quantum chemistry”, ICML 2017.

# Design paradigms for STGNNs

---

Depending on the implementation of the STMP blocks, we categorize STGNNs into:

- **Time-and-Space (T&S)**  
Temporal and spatial processing cannot be factorized in two separate steps.
- **Time-then-Space (TTS)**  
Each time series is embedded in a vector and then representations are propagated on the graph.
- **Space-then-Time (STT)**  
Spatial propagation is performed before processing the resulting time series.

# Time-and-Space

---

In T&S models, representations at every node and time step are obtained by **jointly** propagating representation through time and space.

$$\mathbf{H}_{t-1}^{l+1} = \text{STMP}^l \left( \mathbf{H}_{\leq t-1}^l, \mathcal{E}_{\leq t-1} \right)$$

Several options exist.

- Integrate MP into neural operators for sequential data.
  - Graph recurrent architectures, spatiotemporal convolutions, spatiotemporal attention, ...
- Use sequence molding operators to compute messages.
  - Temporal graph convolutions, spatiotemporal cross-attention, ...
- Product graph representations.

## Example 1: From Recurrent Neural Networks...

Consider a [standard GRU cell](#) [7].

$$\mathbf{r}_t^i = \sigma \left( \Theta_r \left[ \mathbf{x}_t^i || \mathbf{h}_{t-1}^i \right] + \mathbf{b}_r \right) \quad (4)$$

$$\mathbf{u}_t^i = \sigma \left( \Theta_u \left[ \mathbf{x}_t^i || \mathbf{h}_{t-1}^i \right] + \mathbf{b}_u \right) \quad (5)$$

$$\mathbf{c}_t^i = \tanh \left( \Theta_c \left[ \mathbf{x}_t^i || \mathbf{r}_t^i \odot \mathbf{h}_{t-1}^i \right] + \mathbf{b}_c \right) \quad (6)$$

$$\mathbf{h}_t^i = \left( 1 - \mathbf{u}_t^i \right) \odot \mathbf{c}_t^i + \mathbf{u}_t^i \odot \mathbf{h}_{t-1}^i \quad (7)$$

Time series are processed [independently](#) for each node or as a [single multivariate](#) time series.

---

[7] Chung *et al.*, “Empirical evaluation of gated recurrent neural networks on sequence modeling” 2014.

## ...to Graph Convolutional Recurrent Neural Networks

We can obtain a T&S model by implementing the gates of the GRU with MP blocks:

$$\mathbf{Z}_t^l = \mathbf{H}_t^{l-1} \quad (8)$$

$$\mathbf{R}_t^l = \sigma \left( \text{MP}_r^l \left( \left[ \mathbf{Z}_t^l \parallel \mathbf{H}_{t-1}^l \right], \mathcal{E}_t \right) \right), \quad (9)$$

$$\mathbf{O}_t^l = \sigma \left( \text{MP}_o^l \left( \left[ \mathbf{Z}_t^l \parallel \mathbf{H}_{t-1}^l \right], \mathcal{E}_t \right) \right), \quad (10)$$

$$\mathbf{C}_t^l = \tanh \left( \text{MP}_c^l \left( \left[ \mathbf{Z}_t^l \parallel \mathbf{R}_t^l \odot \mathbf{H}_{t-1}^l \right], \mathcal{E}_t \right) \right), \quad (11)$$

$$\mathbf{H}_t^l = \mathbf{O}_t^l \odot \mathbf{H}_{t-1}^l + (1 - \mathbf{O}_t^l) \odot \mathbf{C}_t^l, \quad (12)$$

These T&S models are known as [graph convolutional recurrent neural networks \(GCRNNs\)](#) [8].

---

[8] Seo *et al.*, “Structured sequence modeling with graph convolutional recurrent networks”, ICONIP 2018.



## Popular GCRNNs

---

The **first GCRNN** has been introduced in [8], with **message passing (MP) blocks** implemented as polynomial graph convolutional filters.

GCRNNs have become popular in traffic forecasting with the **Diffusion Convolutional Recurrent Neural Network (DCRNN)** architecture [9].

DCRNN relies on a **bidirectional diffusion convolution**:

$$\mathbf{H}'_t = \sum_{k=0}^K (\mathbf{D}_{t,\text{out}}^{-1} \mathbf{A}_t)^k \mathbf{H}_t \Theta_1^{(k)} + \left( \mathbf{D}_{t,\text{in}}^{-1} \mathbf{A}_t^\top \right)^k \mathbf{H}_t \Theta_2^{(k)} \quad (13)$$

---

[8] Seo *et al.*, “Structured sequence modeling with graph convolutional recurrent networks”, ICONIP 2018.

[9] Li *et al.*, “Diffusion Convolutional Recurrent Neural Network: Data-Driven Traffic Forecasting”, ICLR 2018.

## Example 2: Spatiotemporal convolutional networks (i)

---

Spatiotemporal convolutional networks (STCNs) instead **alternate spatial and temporal convolutions**:

1. Compute intermediate representations by using a **temporal convolutional** layer:

$$z_{t-W:t}^{i,l} = \text{TCN}^l \left( h_{t-W:t}^{i,l-1} \right) \quad \forall i$$

where  $\text{TCN}^l$  indicates a temporal convolutional layer.

2. Compute the updated representation at each time step by using a **graph convolution**:

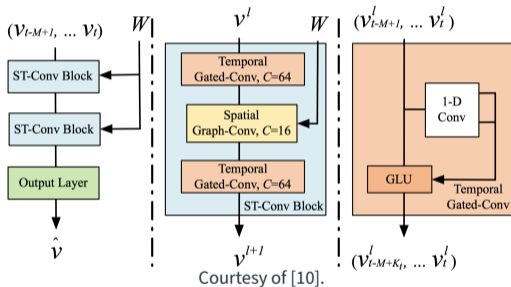
$$H_t^l = \text{MP}^l \left( Z_t^l, \mathcal{E}_t \right) \quad \forall t$$

## Spatiotemporal convolutional networks (ii)

The first example of such architecture is the [STGCN](#) by Yu et al. [10].

The model is obtained by stacking STMP blocks consisting of

- a (gated) temporal convolution;
- a polynomial graph convolution;
- a second (gated) temporal convolution.



More advanced implementations exist, e.g., see [Graph Wavenet](#) [11].

[10] Yu et al., “Spatio-temporal graph convolutional networks: a deep learning framework for traffic forecasting”, IJCAI 2018.

[11] Wu et al., “Graph wavenet for deep spatial-temporal graph modeling”, IJCAI 2019.

## Example 3: Temporal Graph Convolution

A more integrated approach instead consists of implementing a **temporal propagation** mechanism in the message function.

For example, we can design STMP layers s.t.

$$\mathbf{h}_{t-W:t}^{i,l} = \text{TCN}_1^l \left( \mathbf{h}_{t-W:t}^{i,l-1}, \text{AGGR}_{j \in \mathcal{N}_t(i)} \left\{ \text{TCN}_2^l \left( \mathbf{h}_{t-W:t}^{i,l-1}, \mathbf{h}_{t-W:t}^{j,l-1}, \mathbf{e}_{t-W:t}^{ji} \right) \right\} \right).$$

💡 Analogous models can be built with **any sequence modeling architecture**.

→ **Example:** many rely on attention-based operators [12][13].

---

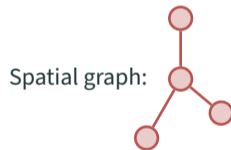
[12] Marisca *et al.*, “Learning to Reconstruct Missing Data from Spatiotemporal Graphs with Sparse Observations”, NeurIPS 2022.

[13] Wu *et al.*, “TraverseNet: Unifying Space and Time in Message Passing for Traffic Forecasting”, TNNLS 2022.

## Example 4: Product graph representations

An alternative option is to consider the sequence  $\mathcal{G}_{t-W:t}$  as a **single graph** with **temporal** and **spatial** edges.

In particular, **product graph representations** can be obtained by **combining the two edge sets**.



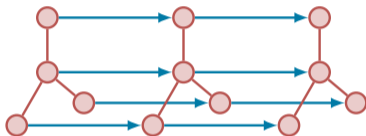
The resulting graph can be processed by any MP neural network.

[14] Sabbaqi *et al.*, “Graph-time convolutional neural networks: Architecture and theoretical analysis”, TPAMI 2023.

# Building product graph representations

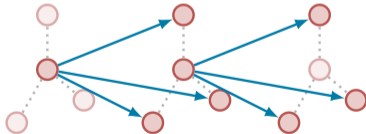
- **Cartesian product**

Spatial graphs are kept and each node is connected to itself in the previous time instant.



- **Kronecker product**

Each node is connected **only** to its neighbors in the previous time instant.

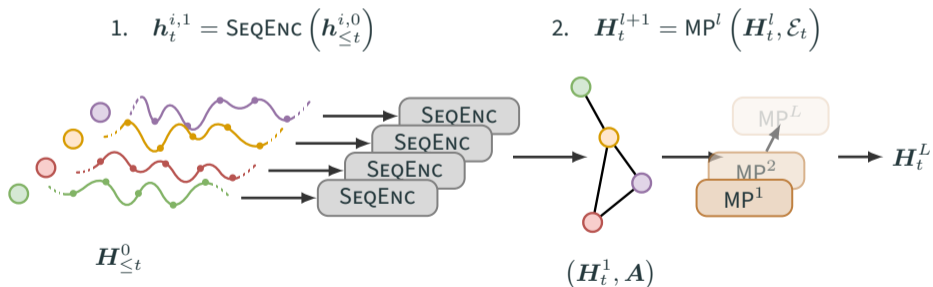


- ...

# Time-then-Space models

The general recipe for a TTS model consists in:

1. **Embedding** each node-level time series in a vector.
2. **Propagating** obtained encodings throughout the graph with a stack of MP layers.



# Pros & Cons of TTS models

---

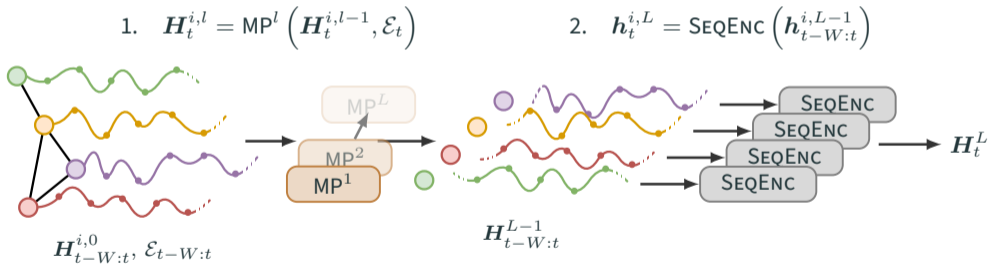
- Pros:**
- 😊 Easy to implement and **computationally efficient**.
  - 😊 We can **reuse operators** we already know.
- Cons:**
- 😞 The 2-step encoding might introduce **information bottlenecks**.
  - 😞 Accounting for **changes in topology** and **dynamic edge attributes** can be more problematic.



# Space-then-Time

In STT approaches the two processing steps of TTS models are inverted:

1. Observations are **propagated among nodes** w.r.t. each time step using a stack of MP layers.
2. Each sequence of representations is processed by a **sequence encoder**.

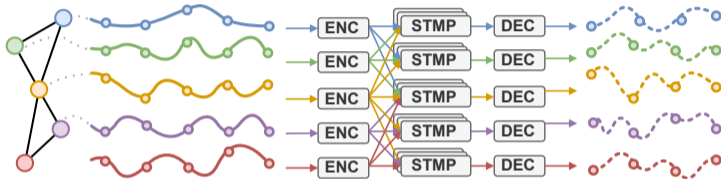


☹️ They do not have the same computational advantages of TTS models.

## **Global and local models**

---

# Globality and locality in STGNNs

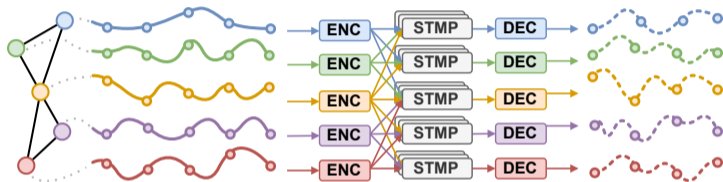


Standard STGNNs are **global** models.

- 😊 Can handle arbitrary node sets.
- 😊 Neighbors provide further conditioning on the predictions.
- 😞 Might struggle with local effects.
- 😞 Might need **long windows** and **high model capacity**.

💡 Use hybrid **global-local** STGNNs.

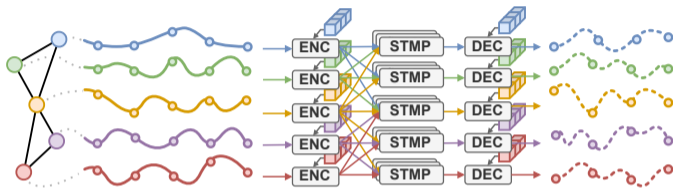
# Global-local STGNNs



💡 We can turn some global components of the architecture into local.

- 😊 Resulting models can capture local effects.
- 😞 Might require a large number of local parameters.

# Global-local STGNNs with node embeddings



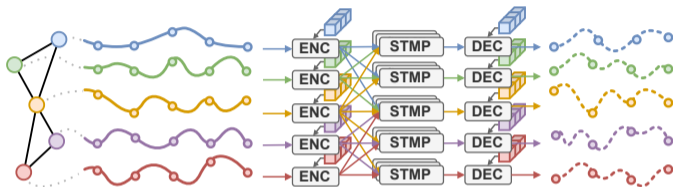
Node embeddings can amortize the learning of local components.

Node embeddings are a table of **learnable parameters**  $Q \in \mathbb{R}^{N \times d_q}$  associated with **each node**.

- 😊 Fed into encoder/decoder, amortize the learning of local components.
- 😊 Most of the model's parameters remain shared.
- 😞 Number of parameters scales linearly with the number of time series . . .
  - One might consider intermediate solutions, e.g., learning embeddings for clusters of time series.

[15] Cini *et al.*, “Taming Local Effects in Graph-based Spatiotemporal Forecasting”, NeurIPS 2023.

# Transferability



! Hybrid global-local STGNNs are not inductive models.

However, the cost of transfer learning can be reduced.

- 😊 Keep shared parameters fixed and finetune local parameters only.
- 😊 Node embeddings can be regularized to facilitate transfer further.

[15] Cini, Marisca, Zambon, and Alippi, “Taming Local Effects in Graph-based Spatiotemporal Forecasting”, NeurIPS 2023.

[16] Butera, De Felice, Cini, and Alippi, “On the Regularization of Learnable Embeddings for Time Series Processing”, Preprint 2024.

## Some empirical results

MODELS	MetrLA	PemsBAY	CER-E	AQI	MetrLA	PemsBAY	CER-E	AQI
Reference arch.	Global models				+ Local node embeddings			
RNN	3.54 $\pm$ .00	1.77 $\pm$ .00	4.57 $\pm$ .01	14.02 $\pm$ .04	<b>3.15</b> $\pm$ .03	<b>1.59</b> $\pm$ .00	<b>4.22</b> $\pm$ .02	<b>13.73</b> $\pm$ .04
GCRNN-IMP	3.35 $\pm$ .01	1.70 $\pm$ .01	4.44 $\pm$ .01	12.87 $\pm$ .02	<b>3.10</b> $\pm$ .01	<b>1.59</b> $\pm$ .00	<b>4.18</b> $\pm$ .01	<b>12.48</b> $\pm$ .03
RNN+IMP	3.34 $\pm$ .01	1.72 $\pm$ .00	4.39 $\pm$ .01	12.74 $\pm$ .02	<b>3.08</b> $\pm$ .01	<b>1.58</b> $\pm$ .00	<b>4.12</b> $\pm$ .03	<b>12.33</b> $\pm$ .02
GCRNN-AMP	3.22 $\pm$ .02	1.65 $\pm$ .00	4.57 $\pm$ .04	12.29 $\pm$ .02	<b>3.07</b> $\pm$ .02	<b>1.59</b> $\pm$ .00	<b>4.17</b> $\pm$ .02	<b>12.17</b> $\pm$ .05
RNN+AMP	3.24 $\pm$ .01	1.66 $\pm$ .00	4.31 $\pm$ .01	12.30 $\pm$ .02	<b>3.06</b> $\pm$ .01	<b>1.58</b> $\pm$ .01	<b>4.13</b> $\pm$ .01	<b>12.15</b> $\pm$ .02
Baseline arch.	Original				+ Local node embeddings			
DCRNN	3.22 $\pm$ .01	1.64 $\pm$ .00	4.28 $\pm$ .01	12.96 $\pm$ .03	<b>3.07</b> $\pm$ .02	<b>1.60</b> $\pm$ .00	<b>4.13</b> $\pm$ .02	<b>12.53</b> $\pm$ .02
GraphWaveNet	3.05 $\pm$ .03	<b>1.56</b> $\pm$ .01	<b>3.97</b> $\pm$ .01	12.08 $\pm$ .11	<b>2.99</b> $\pm$ .02	1.58 $\pm$ .00	4.01 $\pm$ .01	<b>11.81</b> $\pm$ .04
AGCRN	3.16 $\pm$ .01	<b>1.61</b> $\pm$ .00	4.45 $\pm$ .01	13.33 $\pm$ .02	<b>3.14</b> $\pm$ .00	1.62 $\pm$ .00	<b>4.37</b> $\pm$ .02	<b>13.28</b> $\pm$ .03

**Table 1:** MAE on benchmark datasets.

# Transfer learning results

We consider datasets coming from [four different traffic networks](#).

→ **One of the networks is left out** at training time and used for evaluating **transferability**.

RNN+IMP		PEMS03	PEMS04	PEMS07	PEMS08
Fine-tuning	Global	15.30 $\pm$ 0.03	21.59 $\pm$ 0.11	23.82 $\pm$ 0.03	15.90 $\pm$ 0.07
	Embeddings	14.64 $\pm$ 0.05	20.27 $\pm$ 0.11	<b>22.23</b> $\pm$ 0.08	<b>15.45</b> $\pm$ 0.06
	- Variational	<b>14.56</b> $\pm$ 0.03	20.19 $\pm$ 0.05	22.43 $\pm$ 0.02	<b>15.41</b> $\pm$ 0.06
	- Clustering	<b>14.60</b> $\pm$ 0.02	<b>19.91</b> $\pm$ 0.11	<b>22.16</b> $\pm$ 0.07	<b>15.41</b> $\pm$ 0.06
Zero-shot		18.20 $\pm$ 0.09	23.88 $\pm$ 0.08	32.76 $\pm$ 0.69	20.41 $\pm$ 0.07

**Table 2:** Transfer learning results (MAE) after fine-tuning on a week of data.

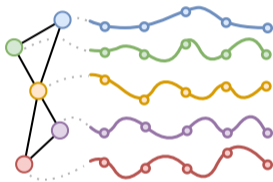
[15] Cini *et al.*, “Taming Local Effects in Graph-based Spatiotemporal Forecasting”, NeurIPS 2023.



# End of Part 1: what we have so far

---

1. We formalized the problem of processing **correlated time series**.
2. **Graph representations** allows for modeling dependencies among them.
3. We discussed **forecasting problem** and **global/local** deep learning models for time series.
4. We saw approaches to building **spatiotemporal graph neural networks** and the associated **trade-offs**.



Before discussing **challenges**, we will look at **software implementations** of the above.

DEMO

# Coding Spatiotemporal GNNs

# tsl: PyTorch Spatiotemporal Library

---



tsl (Torch Spatiotemporal) is a python library built upon [PyTorch](#) and [PyG](#) to accelerate research on neural spatiotemporal data processing methods, with a focus on **Graph Neural Networks**.

 [torch-spatiotemporal.readthedocs.io](https://torch-spatiotemporal.readthedocs.io)

 [github.com/TorchSpatiotemporal/tsl](https://github.com/TorchSpatiotemporal/tsl)



Notebook

**Spatiotemporal Graph Neural Networks with tsl**



---

[17] Cini and Marisca, *Torch Spatiotemporal*, <https://github.com/TorchSpatiotemporal/tsl> 2022.

Part 2

# Challenges

# Challenges

---

- **Scalability**  
How to deal with large collections of time series?
- **Dealing with missing data**  
How to deal with missing observations within the time series?
- **Latent graph learning**  
What to do when the underlying graph is not known?
- **Model quality assessment**  
How to evaluate our graph-based model?

# Scalability

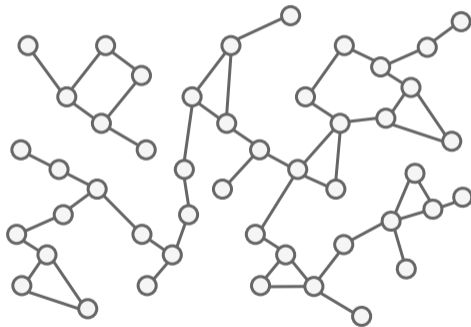
---

## 😊 The scalability feature

---

**Graph-based processing** allows us to

- 😊 learn a single inductive (**global**) model...
- 😊 ...while conditioning on related time series in a **sparse** fashion.
- 😊 The cost of this operation reduces from  $\mathcal{O}(N^2)$  to  $\mathcal{O}(|\mathcal{E}_t|)$



## The scalability issue



---

**Spatiotemporal** data span – as the name suggests – **two dimensions**:

- the **spatial** dimension – the number of time series.
- the **time** dimension – the number of time steps per time series.

In the real world, dealing with **high-frequency, large-scale time series data** is quite common.

– E.g., smart cities, environmental monitoring, finance

-  A large amount of data needs to be **processed at once**.
-  In particular, to account for **long-range spatiotemporal dependencies**.



# Computational complexity of STGNNs

$W$ : length of time series ·  $N$ : number of nodes ·  $|\mathcal{E}_t|$ : number of edges ·  $L$ : number of MP layers

The computational complexity of T&S models is given by:

- node-wise temporal processing –  $\mathcal{O}(WN)$ ;  $\rightarrow \mathcal{O}(W(N + L|\mathcal{E}_t|))$
- $L$  MP layers **for each time step** –  $\mathcal{O}(WL|\mathcal{E}_t|)$ .

A first step toward improving scalability is represented by TTS models, which perform:

- node-wise temporal processing –  $\mathcal{O}(WN)$ ;  $\rightarrow \mathcal{O}(WN + L|\mathcal{E}_t|)$
- $L$  MP layers **at the last time step** –  $\mathcal{O}(L|\mathcal{E}_t|)$ .

STT models, instead, do not have computational advantages over T&S models.

# Graph subsampling

Computations can be reduced by training on **subgraphs** of the full network.

- sampling the  **$K$ -th order neighborhood** of a subset of nodes;
- **rewiring** the graph to reduce the total number of edges.



Mostly adapted from methods developed in **static graph processing** (e.g., [19], [20]).

- ☹️ Subsampling might break long-range spatiotemporal dependencies.
- ☹️ The learning signal may be noisy.

[18] Gandhi *et al.*, “Spatio-Temporal Multi-graph Networks for Demand Forecasting in Online Marketplaces”, ECML-PKDD 2021.

[19] Hamilton *et al.*, “Inductive representation learning on large graphs”, NeurIPS 2017.

[20] Rong *et al.*, “DropEdge: Towards Deep Graph Convolutional Networks on Node Classification”, ICLR 2020.

# Pre-computation

---

Pre-processing methods (e.g., [21]) enable scalability to large graphs by:

- **precomputing** a representation for each **node's neighborhood ahead of training**;
- processing the obtained node representations as if they were **i.i.d. samples**.

An extension to spatiotemporal data is given by **SGP** [22], which acts in 2 steps:

1. obtain a temporal encoding at each time step with a deep **echo state network**<sup>1</sup>;
2. propagate such encodings through the graph using powers of a **graph shift operator**.

---

[21] Frasca *et al.*, “SIGN: Scalable inception graph neural networks” 2020.

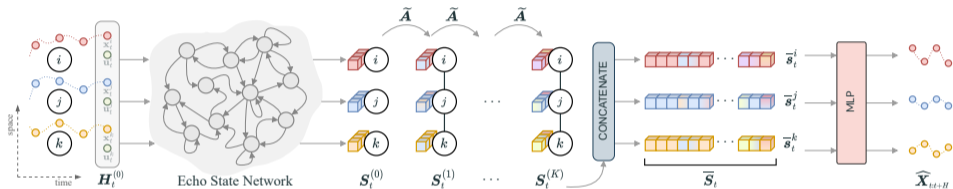
[22] Cini *et al.*, “Scalable Spatiotemporal Graph Neural Networks”, AAAI 2023.

[23] Liu *et al.*, “Do we really need graph neural networks for traffic forecasting?” Preprint 2023.

<sup>1</sup>A randomized recurrent neural networks

# SGP: Scalable Graph Predictor [22]

Extracted representations can be sampled uniformly across time and space during training.



- 😊 The cost of a training step is independent of  $W$ ,  $N$  and  $|\mathcal{E}_t|$ .
- 😊 Performance matches state of the art.
- 😞 More storage space is required – the number of extracted features is  $\gg d_x$ .
- 😞 More reliant on hyperparameter selection than end-to-end approaches.

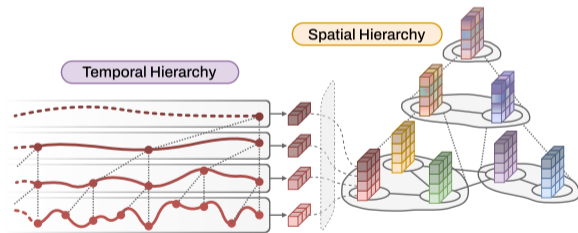
[22] Cini *et al.*, “Scalable Spatiotemporal Graph Neural Networks”, AAI 2023.

# Hierarchical processing

We can reduce computational complexity by using **coarser-grained representations** of the input.

In space, this can be achieved through **graph pooling** [24].

- 😊 Reduced number of operations to reach the same receptive field.
- 😞 Introduce bottlenecks in information propagation.



[24] Grattarola *et al.*, "Understanding Pooling in Graph Neural Networks" 2024.

# Dealing with missing data

---

# The problem of missing data

---

So far, we assumed to deal with **complete sequences**.

- i.e., to have valid observations associated with each node (sensor) and time step.

However, time series collected by real-world sensor networks often have **missing data**, due to:

- faults, of either transient or permanent nature;
- asynchronicity among the time series;
- communication errors...

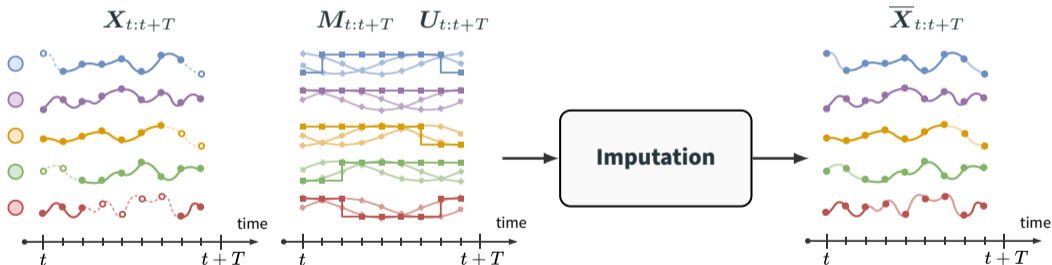
Most forecasting methods operate on complete sequences.

→ We need a way to **impute**, i.e., *reconstruct*, missing data.

# Time series imputation

## Time series imputation (TSI)

Given a window of observations  $\mathbf{X}_{t:t+T}$ , mask  $\mathbf{M}_{t:t+T}$ , and covariates  $\mathbf{U}_{t:t+T}$ , the goal is to estimate the missing observations in the sequence  $\bar{\mathbf{X}}_{t:t+T}$ .



→ We use a **mask**  $m_t^i \in \{0, 1\}$  to distinguish between missing (0) and valid (1) observations.



# Missing data types

We can categorize missing data patterns according to the **conditional distribution**  $p(\mathbf{m}_t^i | \mathbf{M}_{\leq t})$ .

- **Point missing**

$p(\mathbf{m}_t^i = 0)$  is **the same** across nodes and time steps, i.e., RVs associated to each  $\mathbf{m}_t^i$  are iid.

$$p(\mathbf{m}_t^i) = \mathcal{B}(\eta) \quad \forall i, t$$

- **Block missing**

$p(\mathbf{m}_t^i = 0)$  is not independent from missing data **at other nodes and/or time steps**.

**Temporal** block missing  $p(\mathbf{m}_t^i | \mathbf{m}_{t-1}^i) \neq p(\mathbf{m}_t^i)$

**Spatial** block missing  $p(\mathbf{m}_t^i | \{\mathbf{m}_t^j\}^{j \neq i}) \neq p(\mathbf{m}_t^i)$

**Spatiotemporal** block missing  $p(\mathbf{m}_t^i | \mathbf{m}_{t-1}^i, \{\mathbf{m}_t^j\}^{j \neq i}) \neq p(\mathbf{m}_t^i)$

# Optimization

---

Parameters  $\theta$  can be learned by **minimizing a loss function**  $\ell(\cdot, \cdot)$  on **valid observations** in a training set:

$$\hat{\theta} = \arg \min_{\theta} \sum_{t=1}^T \sum_{i=1}^N \frac{\|\mathbf{m}_t^i \odot \ell(\hat{\mathbf{x}}_t^i, \mathbf{x}_t^i)\|_1}{\|\mathbf{m}_t^i\|_1}. \quad \leftarrow \quad \text{e.g., } \ell = (\hat{\mathbf{x}}_t^i - \mathbf{x}_t^i)^2$$

For imputation, we **mark** some valid observations **as missing** with mask  $\overline{\mathbf{m}}_t^i$  to obtain ground-truth labels:

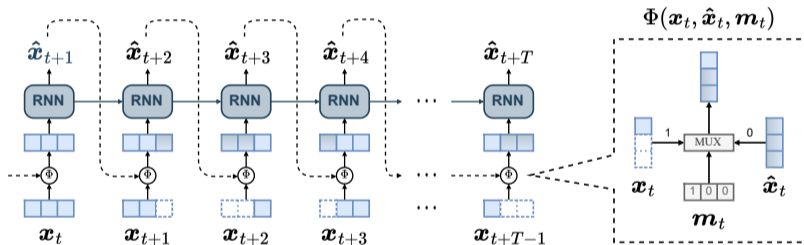
$$\hat{\theta} = \arg \min_{\theta} \sum_{t=1}^T \sum_{i=1}^N \frac{\|\overline{\mathbf{m}}_t^i \odot \ell(\overline{\mathbf{x}}_t^i, \mathbf{x}_t^i)\|_1}{\|\overline{\mathbf{m}}_t^i\|_1}.$$

 Data where  $\overline{\mathbf{m}}_t^i = \mathbf{1}$  must **not be used** in the model to obtain the imputations.

# Deep learning for TSI

Besides standard statistical methods, deep learning approaches have become a popular alternative.

- In particular, **autoregressive models** (e.g., RNNs).

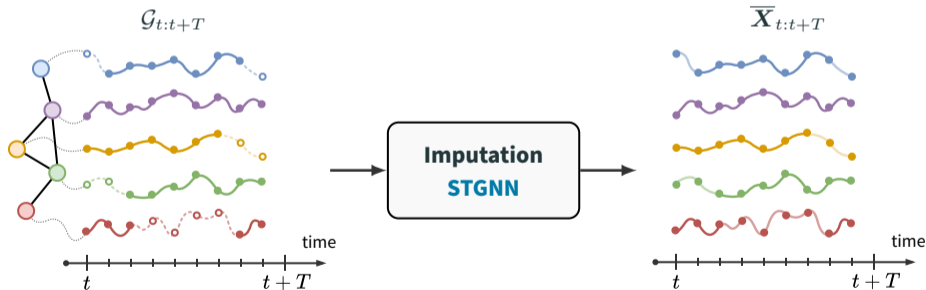


- 😊 Effective in exploiting past (and future, with bidirectional models) **node** observations.
- 😞 Struggle in capturing **nonlinear space-time dependencies**.

# Time series imputation + relational inductive biases

Again, we can use the available relational information to condition the model, i.e.,

$$\mathbf{x}_{t+k}^i \sim p\left(\mathbf{x}_{t+k}^i \mid \mathbf{X}_{t:t+T} \odot \mathbf{M}_{t:t+T}, \mathbf{A}\right) \quad k \in [0, T)$$



# Graph Recurrent Imputation Network (GRIN)

Similarly to GCRNN for forecasting, we can integrate graph processing into the autoregressive approach for imputation [25].

In these approaches, the distribution  $p(\mathbf{x}_t^i | \mathbf{X}_{0:\infty} \odot \mathbf{M}_{0:\infty})$  is modeled into **three independent steps**:

Information from  
previous observations.

$$p(\mathbf{x}_t^i | \mathbf{X}_{<t} \odot \mathbf{M}_{<t})$$

Typically modeled by bidirectional autoregressive models.

Information from  
subsequent observations.

$$p(\mathbf{x}_t^i | \mathbf{X}_{>t} \odot \mathbf{M}_{>t})$$

Information from related  
concurrent observations.

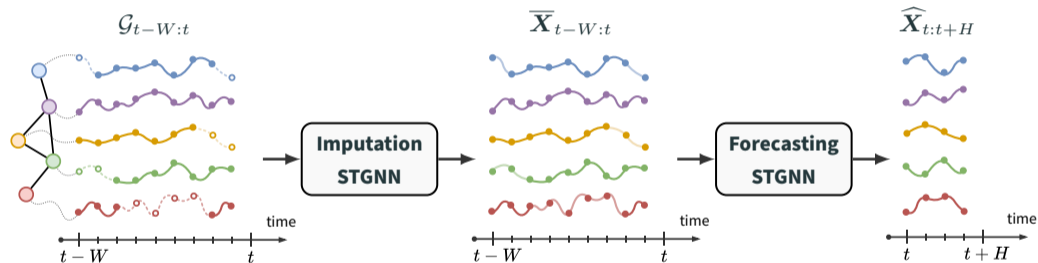
$$p(\mathbf{x}_t^i | \{\mathbf{x}_t^j \odot \mathbf{m}_t^j\}^{j \neq i})$$

Enabled by message passing.

[25] Cini *et al.*, "Filling the G\_ap\_s: Multivariate Time Series Imputation by Graph Neural Networks", ICLR 2022.

## Imputation *before* forecasting

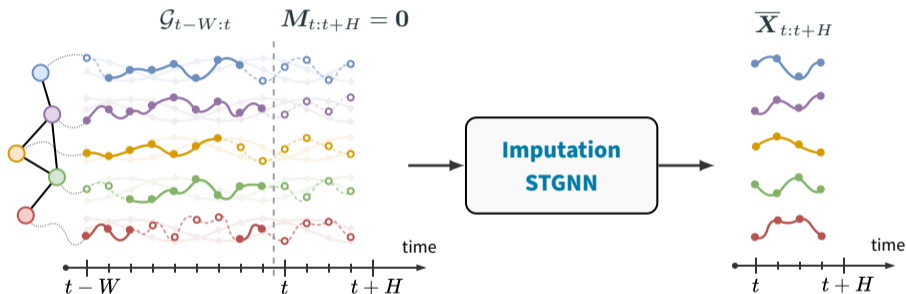
TSI is often used as a **preprocessing step** for a downstream task, e.g., forecasting.



- ☹️ Often necessary to use standard forecasting methods with irregular time series.
- ☹️ Might introduce **biases** due to errors in estimated values.

## Imputation *in place of* forecasting

Imputation methods can also be adapted to perform forecasting.



- ☹ It is a **workaround** (this is not their purpose).
- ☹ Might perform poorly due to the absence of values in the forecasting horizon.

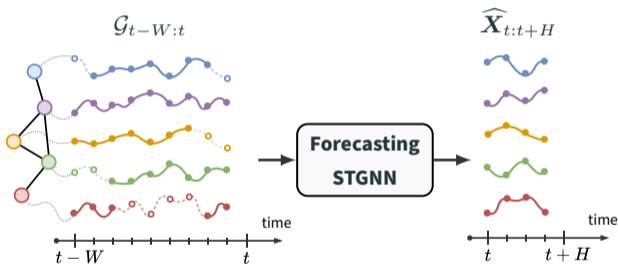
# Forecasting from partial observations

A more direct approach: **avoid the reconstruction step!**

→ Design forecasting architecture to **directly deal with irregular observations.**

## Benefits

- 😊 Learn how to leverage **only valid observations** specifically for the task at hand.
- 😊 Avoid the computational burden of **imputing missing values.**



[26] Zhang *et al.*, “Graph-guided network for irregularly sampled multivariate time series”, ICLR 2022.

[27] Zhong *et al.*, “Heterogeneous spatio-temporal graph convolution network for traffic forecasting with missing values”, IEEE ICDCS 2021.

[28] Marisca *et al.*, “Graph-based Forecasting with Missing Data through Spatiotemporal Downsampling”, ICML 2024.

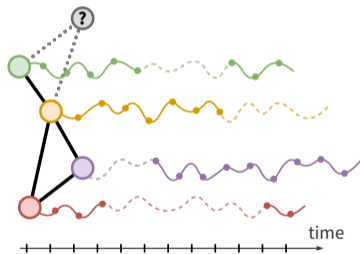


# Virtual sensing

The practice of estimating unmeasured states using models and existing observations.

The power of graphs:

- 😊 The **relational** processing allows us to condition estimates on data **close in space**.
- 😊 The **inductive** property of MP allows us to handle **new nodes and edges**.
- 😊 Useful in applications where sensing has a cost.



[29] Wu *et al.*, “Inductive Graph Neural Networks for Spatiotemporal Kriging”, AAAI 2021.

[30] De Felice *et al.*, “Graph-Based Virtual Sensing from Sparse and Partial Multivariate Observations”, ICLR 2024.

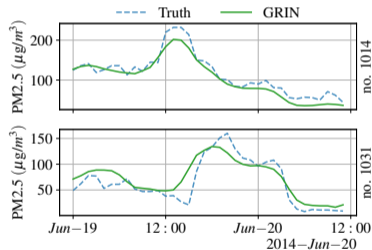
# Graph imputation for virtual sensing

💡 Add a **fictitious node** with **no data** and let the model **infer** the corresponding time series.

Clearly, several assumptions are needed

- high degree of homogeneity of sensors,
- capability to reconstruct from observations at neighboring sensors,
- and many more...

Two virtual sensors for air quality. (from [25])



[12] Marisca *et al.*, “Learning to Reconstruct Missing Data from Spatiotemporal Graphs with Sparse Observations”, NeurIPS 2022.

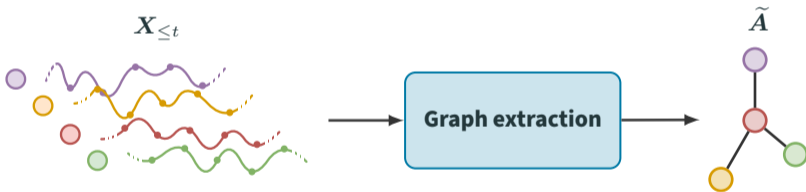
[25] Cini *et al.*, “Filling the G<sub>ap</sub>s: Multivariate Time Series Imputation by Graph Neural Networks”, ICLR 2022.

# Latent graph learning

---

# Learning an adjacency matrix

- ☹ Relational information is **not** always (or only partially) **available**,
- ☹ or might be **ineffective** in capturing spatial dynamics.
- 😊 **Relational** architectural **biases** can nonetheless be exploited  
→ **extract a graph** from the time series or node attributes



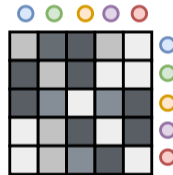
- When possible, the learned graph should be **sparse**.
- It can be interpreted as **regularizing** a **spatial attention** operator.
- This task is found under different names:  
**graph structure learning**, latent graph learning, graph inference...

# Time-series similarities

---

Probably, the simplest approach to extract a graph from the time series is by computing **time series similarity scores**.

- Pearson correlation
- Correntropy
- Granger causality
- Kernels for time series
- ...



→ Thresholding might be necessary to obtain binary and sparse graphs.

# Inferring latent structures from time series

---

Model the **graph as a latent variable** determining the realizations of the time series.

- They rely on assumptions, such as of signal smoothness and of a diffusion process.

Dedicated **loss functions** are formulated and minimized, e.g.,

$$\text{trace}(\mathbf{X}^\top \mathbf{L} \mathbf{X}) = \frac{1}{2} \sum_{ij} \mathbf{A}_{i,j} \|\mathbf{X}_i - \mathbf{X}_j\|_2^2$$

constraining  $\mathbf{L}$  (or  $\mathbf{A}$ ) to be a Laplacian (adjacency matrix) and promoting sparsity.

→ These approaches are commonly derived from a graph signal processing point of view.

---

[31] Dong *et al.*, “Learning Laplacian matrix in smooth graph signal representations”, IEEE TSP 2016.

[32] Mateos *et al.*, “Connecting the dots: Identifying network structure via graph signal processing”, IEEE SP Mag 2019.

# Task-oriented latent graph learning


---

An integrated approach: **learn** the **relations** **end-to-end** with the downstream task

→ e.g., by minimizing the forecasting error (MAE, MSE...).

Two different formulations:

1. learning directly an **adjacency matrix**  $\mathbf{A} \in \mathbb{R}^{N \times N}$ ;
2. learning a **probability distribution over graphs**  $p_{\Phi}$  generating  $\mathbf{A}$  (often  $\in \{0, 1\}^{N \times N}$ ).

 One key challenge is keeping both  $\mathbf{A}$  and the subsequent computations **sparse**.  
→ **non-trivial** with gradient-based optimization.

# Direct approach

A direct approach consists in learning  $\tilde{\mathbf{A}}$  as function  $\xi(\cdot)$  of edge scores

$\Phi \in \mathbb{R}^{N \times N}$  as

$$\tilde{\mathbf{A}} = \xi(\Phi)$$

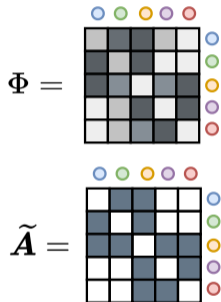
Edge scores  $\Phi$

- can be a table of **learnable** model **parameters**,
- obtained as a **function** of the **inputs** and/or other parameters:

$$\Phi = \Phi(\mathbf{X}, \tilde{\Phi}).$$

Function  $\xi(\cdot)$  can enforce structures on  $\tilde{\mathbf{A}}$ , like,

→ make  $\tilde{\mathbf{A}}$  **binary**, a  **$k$ -NN** graph, a **tree**...





# Edge score factorization

The number of possible edge scores is **quadratic** in the number of nodes ( $\Phi \in \mathbb{R}^{N \times N}$ )

→ a common approach is to factorize  $\Phi$ :

$$\tilde{\mathbf{A}} = \xi(\Phi) \quad \Phi = \mathbf{Z}_s \mathbf{Z}_t^\top$$

with

- $\mathbf{Z}_s \in \mathbb{R}^{N \times d}$  **source** node embeddings
- $\mathbf{Z}_t \in \mathbb{R}^{N \times d}$  **target** node embeddings

$$\Phi = \mathbf{Z}_s \mathbf{Z}_t^\top$$

$\mathbf{Z}_s$  and  $\mathbf{Z}_t$  can be learned as tables of (local) parameters or **as a function of the input window**.

[11] Wu *et al.*, “Graph wavenet for deep spatial-temporal graph modeling”, IJCAI 2019.

## Pro & Cons of the direct approach

---

- 😊 Easy to implement.
- 😊 Many possible parametrizations.
- 😊 Edge scores are usually easy to learn end-to-end.
- 😞 It often results in dense computations with  $\mathcal{O}(N^2)$  complexity.
- 😞 Sparsifying  $A$  results in sparse gradients.
- 😞 Encoding prior structural information requires smart parametrizations.

## Probabilistic methods

In this context, probabilistic methods aim at learning a **parametric distribution**  $p_{\Phi}$  for  $\mathbf{A}$ .

- Different parametrizations of  $p_{\Phi}$  allow for embedding **graph structural priors** on the sampled graphs, e.g., edge density, bound node degrees.

### Graphs of independent edges

For every edge  $(i, j)$

$$\mathbf{A}_{i,j} \sim \text{Bernoulli}(\sigma(\Phi_{i,j})).$$

### Fixed-degree graphs

For each node  $i$ , sample w/o replacement  $k$  nodes from

$$\text{Categorical}(\text{SoftMax}(\Phi_{i,1}, \dots, \Phi_{i,N})).$$

- As seen before,  $\Phi$  can be factorized and  $p_{\Phi}$  made input dependent, e.g.,

$$\Phi = \xi \left( \mathbf{Z}_s \mathbf{Z}_t^{\top} \right), \quad \mathbf{A} \sim p_{\Phi}(\mathbf{A} | \mathbf{X}_{<t}, \mathbf{U}_{<t}, \mathbf{V}).$$

[33] Kazi *et al.*, “Differentiable graph module (dgm) for graph convolutional networks”, IEEE TPAMI 2022.

[34] Cini *et al.*, “Sparse graph learning from spatiotemporal time series”, JMLR 2023.

# Learning graph distributions

Training losses average over all graphs according to  $p_{\Phi}$ , e.g., based on point predictions

$$\mathcal{L}(\Phi) = \mathbb{E}_{\mathbf{A} \sim p_{\Phi}} \left[ \ell \left( \widehat{\mathbf{X}}_{t:t+H}, \mathbf{X}_{t:t+H} \right) \right], \quad \mathcal{L}(\Phi) = \ell \left( \mathbb{E}_{\mathbf{A} \sim p_{\Phi}} \left[ \widehat{\mathbf{X}}_{t:t+H} \right], \mathbf{X}_{t:t+H} \right),$$

where  $\mathbf{X}_{t:t+H} = \mathcal{F}(\mathbf{X}_{t-W:t}, \mathbf{A}; \theta)$ .

More generally, comparing predictive distributions by means of a divergence measure  $\Delta$

$$\mathcal{L}(\Phi) = \Delta \left( p_{\Phi}(\widehat{\mathbf{X}}_{t:t+H}), p(\mathbf{X}_{t:t+H}) \right).$$



Gradient-based optimization requires computing  $\nabla_{\Phi} \mathcal{L}(\Phi)$ ,  
 → i.e., differentiating w.r.t. the parameters of the integrated distribution.

- ☹ Analytical computations is often unfeasible;
- ☹ Monte Carlo approximations require care.

## Monte Carlo gradient estimators

💡 One approach is to **reparametrize**  $\tilde{\mathbf{A}} \sim p_{\Phi}(\mathbf{A})$  as:  $\tilde{\mathbf{A}} = g(\Phi, \epsilon)$ ,  $\epsilon \sim p(\epsilon)$   
 decoupling parameters  $\Phi$  from the random component  $\epsilon$ :  $\nabla_{\Phi} \mathcal{L}(\Phi) = \mathbb{E}_{\epsilon} \left[ \nabla_{\Phi} \ell(\widehat{\mathbf{X}}, \mathbf{X}) \right]$ .

😊 Practical and **easy** to implement,

😞 rely on **continuous relaxations** and make subsequent computations scale with  $\mathcal{O}(N^2)$ .

💡 Conversely, **score-function** (SF) gradient estimators rely on the relation

$$\nabla_{\Phi} \mathbb{E}_{p_{\Phi}} \left[ \ell(\widehat{\mathbf{X}}, \mathbf{X}) \right] = \mathbb{E}_{p_{\Phi}} \left[ \ell(\widehat{\mathbf{X}}, \mathbf{X}) \nabla_{\Phi} \log p_{\Phi} \right]$$

😞 suffer from **high variance** (use variance reduction techniques),

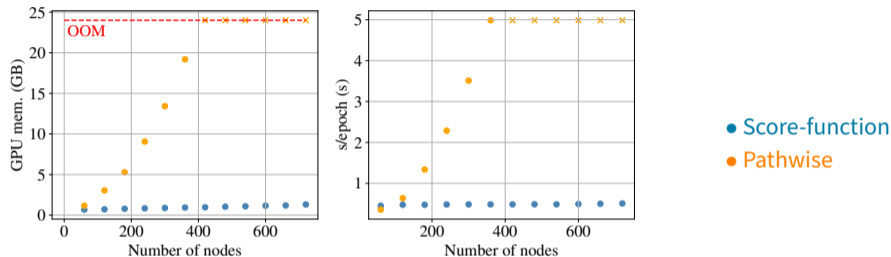
😊 allow to **keep computations sparse** through the model.

→ we can use **Monte Carlo** estimator.

[35] Kipf *et al.*, “Neural relational inference for interacting systems”, ICML 2018.

[34] Cini *et al.*, “Sparse graph learning from spatiotemporal time series”, JMLR 2023.

# Computational efficiency



With score-based gradient estimators  $\nabla_{\Phi} \mathcal{L}(\Phi) = \mathbb{E}_{\mathbf{A} \sim p_{\Phi}} \left[ \ell(\widehat{\mathbf{X}}, \mathbf{X}) \nabla_{\Phi} \log p_{\Phi}(\mathbf{A}) \right]$ .

- ☺ They are **computationally efficient** as  $\nabla_{\Phi}$  is computed with respect to  $\log p_{\Phi}(\mathbf{A})$ .
  - do not rely on continuous relaxation of **discrete random variables**;
  - allow for **sparse message passing** to compute  $\widehat{\mathbf{X}}$  (and, in turn, of  $\ell(\widehat{\mathbf{X}}, \mathbf{X})$ ) by rely on sparse matrices  $\mathbf{A}$ .
- ☹ They can be sample inefficient due to the **high variance** of the gradient estimates.

## Uncertainty quantification

While probabilistic models have been used to enable learning of discrete variables (graph edges), the associated **edge probabilities** can carry information about the **relevance** of the associated connections.

→ It enables some degree of explainability and better informed decision-making.

⚠ Assessing the calibration of latent variables is **hard** on real data.

→ This is due to their latent nature, for which observations are not available.

💡 Studies provide some learning guarantees, e.g.,

Minimizing appropriate divergence measures  $\Delta \left( p(\mathbf{X}), p_{\Phi}(\widehat{\mathbf{X}}) \right)$  of the data and predictive distributions  $p(\mathbf{X}), p_{\Phi}(\widehat{\mathbf{X}})$ , respectively, enables calibration of the  $p_{\Phi}(\mathbf{A})$ .

[36] Gray *et al.*, “Bayesian inference of network structure from information cascades”, IEEE TSIPN 2020.

[37] Manenti *et al.*, *Learning Latent Graph Structures and Their Uncertainty*, Preprint 2024.

# **Model quality assessment**

---



## Questions to answer

---

Consider a predictor  $\mathcal{F}$  trained to solve a time-series forecasting problem.

1. Is the predictor **optimal** for the problem at hand?
2. **Where** does the predictor appear sub-optimal?
3. **How** can we improve the predictor?

**Remark:** Multiple optimality criteria can be considered.

 Relational inductive biases can help us here too.

# Performance at task

---

Given two predictors  $\mathcal{F}_a, \mathcal{F}_b$  and performance metric  $M$  (e.g., MAE, MSE).

- we consider  $\mathcal{F}_a$  **better** than  $\mathcal{F}_b$  if  $M(\mathcal{F}_a)$  is *statistically* better than  $M(\mathcal{F}_b)$ .
- we consider  $\mathcal{F}_a$  **optimal** if there is no other model  $\mathcal{F}_b$  better than  $\mathcal{F}_a$ .

Can we further improve over the best model so far  $\mathcal{F}_a$ ?

- Either we **find a new model**  $\mathcal{F}_*$  better than  $\mathcal{F}_a$
- or we need **prior knowledge** about the modeled system.

Model	$M$
$\mathcal{F}_a$	$0.145_{\pm 0.002}$
$\mathcal{F}_b$	$0.176_{\pm 0.005}$
$\vdots$	
$\mathcal{F}_n$	$0.158_{\pm 0.004}$
$\mathcal{F}_*$	$0.139_{\pm 0.001}$

## Residual correlation analysis

Studying the **correlation** between prediction residuals  $r_t^i \doteq \mathbf{x}_{t:t+H}^i - \hat{\mathbf{x}}_{t:t+H}^i$  allows for testing model optimality.

If residuals are **dependent**

⇒ there is **information** that the model **hasn't captured**

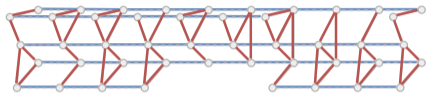
⇒ model predictions **can be improved**.

**Remarks:** Residual correlation analysis

- 😊 Is independent of specific performance measures.
- 😞 Does not quantify how much a model can improve w.r.t. a specific performance metric.
- 😊 Does not rely on comparisons with other models.

Research focused mainly on either serial correlation [38]–[40] or spatial correlation [41], [42].

# AZ-Whiteness test: a spatio-temporal test



The test is defined by statistic

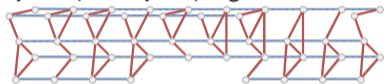
$$C(\{\mathbf{r}\}) = \underbrace{\sum_t \sum_{(i,j) \in \mathcal{E}_t} w_{ijt} \operatorname{sgn}(\langle \mathbf{r}_t^i, \mathbf{r}_t^j \rangle)}_{\text{spatial edge}} + \underbrace{\sum_t \sum_i w_{it} \operatorname{sgn}(\langle \mathbf{r}_t^i, \mathbf{r}_{t+1}^i \rangle)}_{\text{temporal edge}} \rightarrow \mathcal{N}(0, 1)$$

- ☺ distribution-free and residuals can be non-identically distributed.
- ☺ computation is linear in the number of edges and time steps.

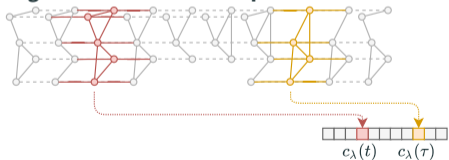
# Where can we improve?

Analyzing the AZ-whiteness test statistic computed on subgraphs of the spatio-temporal graph allows for discovering insightful correlation patterns.

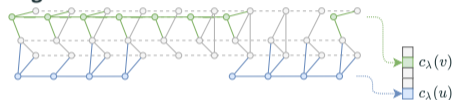
## Spatial (or temporal) edges



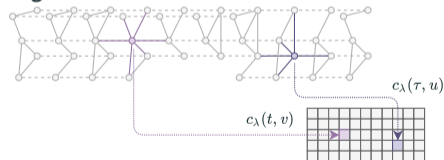
## Edges related to a time step



## Edges related to a node



## Edges related to a node

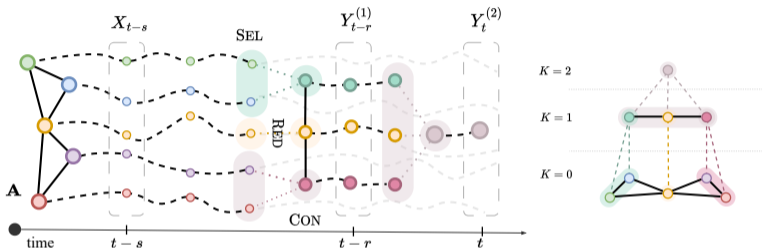


[44] Zambon *et al.*, “Where and How to Improve Graph-based Spatio-temporal Predictors” 2023.

## **Future directions**

---

# Hierarchical processing



- ☹️ Standard STGNNs operate at a **fixed spatiotemporal scale**.
- 💡 Combine **hierarchical** and **graph-based** representations.
- 😊 Exploit **higher-order dependencies** by operating on **hierarchical representations** of the input.
- 😊 Can also be used for **hierarchical forecasting** and to obtain **reconciled predictions**.

[45] Yu *et al.*, “ST-Unet: A spatio-temporal U-network for graph-structured time series modeling” 2019.

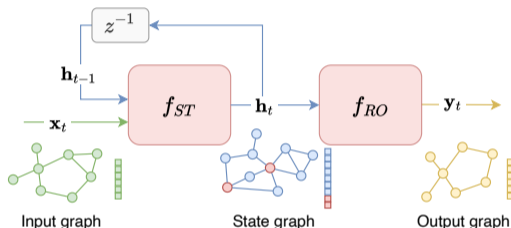
[46] Cini *et al.*, “Graph-based Time Series Clustering for End-to-End Hierarchical Forecasting”, ICML 2024.

[28] Marisca *et al.*, “Graph-based Forecasting with Missing Data through Spatiotemporal Downsampling”, ICML 2024.

# State-space models

$$\begin{cases} \mathbf{h}_t = f_{ST}(\mathbf{h}_{t-1}, \mathbf{x}_{t-1}, \boldsymbol{\eta}_{t-1}) \\ \mathbf{y}_t = f_{RO}(\mathbf{h}_t, \boldsymbol{\nu}_t) \end{cases}$$

- Inputs  $\mathbf{x}_t$ , states  $\mathbf{h}_t$ , and outputs  $\mathbf{y}_t$  are different attributed graphs.
- $\boldsymbol{\eta}_t, \boldsymbol{\nu}_t$  are noise terms at the node/edge level.



[47] Rangapuram *et al.*, “Deep State Space Models for Time Series Forecasting”, NeurIPS 2018.

[48] Zambon *et al.*, *Graph State-Space Models*, Preprint 2023.

[49] Alippi *et al.*, *Graph Kalman Filters*, Preprint 2023.

[50] Buchnik *et al.*, “GSP-kalmanet: Tracking graph signals via neural-aided Kalman filtering”, IEEE TSP 2024.

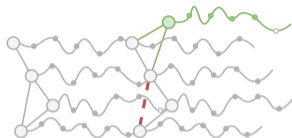
[51] Chouzenoux *et al.*, “Sparse graphical linear dynamical systems”, JMLR 2024.



# Inductive learning

In real-world applications, one often needs to

- operate under **changes** in the network **connectivity**
- make predictions for **newly added nodes**
- **transfer** the model to **different** sensor **networks** (collections of time series)



Useful in **several tasks**, like, forecasting, missing data imputation, and virtual sensing.

**!** **Performance** can easily **degrade** if the **data distribution** of target nodes

- **deviates** from that at **training nodes**
- **changes over time**.

[15] Cini *et al.*, “Taming Local Effects in Graph-based Spatiotemporal Forecasting”, NeurIPS 2023.

[52] Yin *et al.*, “Nodetrans: A graph transfer learning approach for traffic prediction”, Preprint 2022.

[53] Prabowo *et al.*, “Traffic forecasting on new roads using spatial contrastive pre-training (SCPT)” 2024.

# Benchmarks

---

## Open datasets

In line with OGB [54], TGB [55], TGB 2.0 [56].

- Energy analytics (CER-E, PV-US) [22]
- Traffic flow (LargeST) [57]
- ...

## Software

Standard model evaluation platforms

- Torch SpatioTemporal [17]
- BasicTS [58]
- ...

---

[22] Cini *et al.*, “Scalable Spatiotemporal Graph Neural Networks”, AAI 2023.

[57] Liu *et al.*, “Largest: A benchmark dataset for large-scale traffic forecasting”, NeurIPS (D&B) 2024.

[17] Cini *et al.*, *Torch Spatiotemporal*, <https://github.com/TorchSpatiotemporal/tsl> 2022.

[58] Shao *et al.*, “Exploring Progress in Multivariate Time Series Forecasting: Comprehensive Benchmarking and Heterogeneity Analysis”, IEEE TKDE 2024.

# Conclusions

---

## Some Takeaways

---



- ⚡ **Relational** inductive **biases** allow for exploiting dependencies among the time series,
- 😊 ...while **sharing** most of the model **parameters**,
- 😊 ...and overcoming limits due to **irregularities in time and space**.
- 💡 Whenever possible, **global-local models** are a safe starting point.

**Challenges.** Scalability • Missing data • Latent graph learning • Model quality assessment

**Resources.** 📄 Tutorial paper [3] • 🔄 Open-source library [17]

---

[3] Cini, Marisca, Zambon, and Alippi, “Graph Deep Learning for Time Series Forecasting”, Preprint 2023.

[17] Cini and Marisca, *Torch Spatiotemporal*, <https://github.com/TorchSpatiotemporal/tsl> 2022.



**Andrea Cini**

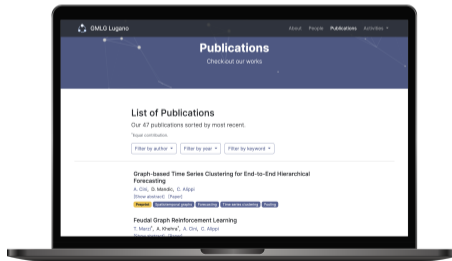


**Ivan Marisca**



**Daniele Zambon**

**Graph Machine Learning Group**  
[gmlg.ch](https://gmlg.ch)



# THE END

Questions?

# References i

---

- [1] D. Salinas, V. Flunkert, J. Gasthaus, and T. Januschowski, “**DeepAR: Probabilistic forecasting with autoregressive recurrent networks,**” *International Journal of Forecasting*, vol. 36, no. 3, pp. 1181–1191, 2020.
- [2] K. Benidis, S. S. Rangapuram, V. Flunkert, *et al.*, “**Deep learning for time series forecasting: Tutorial and literature survey,**” *ACM Comput. Surv.*, vol. 55, no. 6, Dec. 2022, ISSN: 0360-0300. DOI: 10.1145/3533382. [Online]. Available: <https://doi.org/10.1145/3533382>.
- [3] A. Cini, I. Marisca, D. Zambon, and C. Alippi, “**Graph deep learning for time series forecasting,**” *arXiv preprint arXiv:2310.15978*, 2023.
- [4] P. Montero-Manso and R. J. Hyndman, “**Principles and algorithms for forecasting groups of time series: Locality and globality,**” *International Journal of Forecasting*, vol. 37, no. 4, pp. 1632–1653, 2021.
- [5] R. Sen, H.-F. Yu, and I. S. Dhillon, “**Think globally, act locally: A deep neural network approach to high-dimensional time series forecasting,**” *Advances in Neural Information Processing Systems*, vol. 32, 2019.

## References ii

---

- [6] J. Gilmer, S. S. Schoenholz, P. F. Riley, O. Vinyals, and G. E. Dahl, “**Neural message passing for quantum chemistry,**” in *International conference on machine learning*, PMLR, 2017, pp. 1263–1272.
- [7] J. Chung, C. Gulcehre, K. Cho, and Y. Bengio, “**Empirical evaluation of gated recurrent neural networks on sequence modeling,**” *arXiv preprint arXiv:1412.3555*, 2014.
- [8] Y. Seo, M. Defferrard, P. Vandergheynst, and X. Bresson, “**Structured sequence modeling with graph convolutional recurrent networks,**” in *International Conference on Neural Information Processing*, Springer, 2018, pp. 362–373.
- [9] Y. Li, R. Yu, C. Shahabi, and Y. Liu, “**Diffusion convolutional recurrent neural network: Data-driven traffic forecasting,**” in *International Conference on Learning Representations*, 2018.
- [10] B. Yu, H. Yin, and Z. Zhu, “**Spatio-temporal graph convolutional networks: A deep learning framework for traffic forecasting,**” in *Proceedings of the 27th International Joint Conference on Artificial Intelligence*, 2018, pp. 3634–3640.
- [11] Z. Wu, S. Pan, G. Long, J. Jiang, and C. Zhang, “**Graph wavenet for deep spatial-temporal graph modeling,**” in *Proceedings of the 28th International Joint Conference on Artificial Intelligence*, 2019, pp. 1907–1913.



## References iii

---

- [12] I. Marisca, A. Cini, and C. Alippi, “**Learning to reconstruct missing data from spatiotemporal graphs with sparse observations,**” in *Advances in Neural Information Processing Systems*, 2022.
- [13] Z. Wu, D. Zheng, S. Pan, Q. Gan, G. Long, and G. Karypis, “**Traversenet: Unifying space and time in message passing for traffic forecasting,**” *IEEE Transactions on Neural Networks and Learning Systems*, 2022.
- [14] M. Sabbaqi and E. Isufi, “**Graph-time convolutional neural networks: Architecture and theoretical analysis,**” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 45, no. 12, pp. 14 625–14 638, Dec. 2023, ISSN: 1939-3539. DOI: 10.1109/TPAMI.2023.3311912.
- [15] A. Cini, I. Marisca, D. Zambon, and C. Alippi, “**Taming local effects in graph-based spatiotemporal forecasting,**” *arXiv preprint arXiv:2302.04071*, 2023.
- [16] L. Butera, G. De Felice, A. Cini, and C. Alippi, “**On the regularization of learnable embeddings for time series processing,**” *arXiv preprint arXiv:2410.14630*, 2024.
- [17] A. Cini and I. Marisca, *Torch Spatiotemporal*, Mar. 2022. [Online]. Available: <https://github.com/TorchSpatiotemporal/tsl>.

# References iv

---

- [18] A. Gandhi, Aakanksha, S. Kaveri, and V. Chaoji, “**Spatio-temporal multi-graph networks for demand forecasting in online marketplaces,**” in *Machine Learning and Knowledge Discovery in Databases. Applied Data Science Track: European Conference, ECML PKDD 2021, Bilbao, Spain, September 13–17, 2021, Proceedings, Part IV*, 2021, pp. 187–203, ISBN: 978-3-030-86513-9. DOI: 10.1007/978-3-030-86514-6\_12. [Online]. Available: [https://doi.org/10.1007/978-3-030-86514-6\\_12](https://doi.org/10.1007/978-3-030-86514-6_12).
- [19] W. Hamilton, Z. Ying, and J. Leskovec, “**Inductive representation learning on large graphs,**” *Advances in neural information processing systems*, vol. 30, 2017.
- [20] Y. Rong, W. Huang, T. Xu, and J. Huang, “**Dropedge: Towards deep graph convolutional networks on node classification,**” in *International Conference on Learning Representations*, 2020.
- [21] F. Frasca, E. Rossi, D. Eynard, B. Chamberlain, M. Bronstein, and F. Monti, “**SIGN: Scalable inception graph neural networks,**” *arXiv preprint arXiv:2004.11198*, 2020.
- [22] A. Cini, I. Marisca, F. M. Bianchi, and C. Alippi, “**Scalable spatiotemporal graph neural networks,**” *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 37, no. 6, pp. 7218–7226, Jun. 2023. DOI: 10.1609/aaai.v37i6.25880.

# References v

---

- [23] X. Liu, Y. Liang, C. Huang, *et al.*, “**Do we really need graph neural networks for traffic forecasting?**” *arXiv preprint arXiv:2301.12603*, 2023.
- [24] D. Grattarola, D. Zambon, F. M. Bianchi, and C. Alippi, “**Understanding pooling in graph neural networks,**” *IEEE Transactions on Neural Networks and Learning Systems*, vol. 35, no. 2, pp. 2708–2718, 2024.
- [25] A. Cini, I. Marisca, and C. Alippi, “**Filling the g\_ap\_s: Multivariate time series imputation by graph neural networks,**” in *International Conference on Learning Representations*, 2022. [Online]. Available: <https://openreview.net/forum?id=k0u3-S3wJ7>.
- [26] X. Zhang, M. Zeman, T. Tsiligkaridis, and M. Zitnik, “**Graph-guided network for irregularly sampled multivariate time series,**”, 2022.
- [27] W. Zhong, Q. Suo, X. Jia, A. Zhang, and L. Su, “**Heterogeneous spatio-temporal graph convolution network for traffic forecasting with missing values,**” in *2021 IEEE 41st International Conference on Distributed Computing Systems (ICDCS)*, IEEE, 2021, pp. 707–717.
- [28] I. Marisca, C. Alippi, and F. M. Bianchi, “**Graph-based forecasting with missing data through spatiotemporal downsampling,**” in *Proceedings of the 41st International Conference on Machine Learning*, ser. *Proceedings of Machine Learning Research*, vol. 235, PMLR, 2024, pp. 34 846–34 865.

# References vi

---

- [29] Y. Wu, D. Zhuang, A. Labbe, and L. Sun, “**Inductive Graph Neural Networks for Spatiotemporal Kriging,**” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 35, 2021, pp. 4478–4485.
- [30] G. De Felice, A. Cini, D. Zambon, V. Gusev, and C. Alippi, “**Graph-based Virtual Sensing from Sparse and Partial Multivariate Observations,**” in *The Twelfth International Conference on Learning Representations*, 2024. [Online]. Available: <https://openreview.net/forum?id=CAqdG2dy5s>.
- [31] X. Dong, D. Thanou, P. Frossard, and P. Vandergheynst, “**Learning laplacian matrix in smooth graph signal representations,**” *IEEE Transactions on Signal Processing*, vol. 64, no. 23, pp. 6160–6173, 2016.
- [32] G. Mateos, S. Segarra, A. G. Marques, and A. Ribeiro, “**Connecting the dots: Identifying network structure via graph signal processing,**” *IEEE Signal Processing Magazine*, vol. 36, no. 3, pp. 16–43, 2019.
- [33] A. Kazi, L. Cosmo, S.-A. Ahmadi, N. Navab, and M. M. Bronstein, “**Differentiable graph module (dgm) for graph convolutional networks,**” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 45, no. 2, pp. 1606–1617, 2022.
- [34] A. Cini, D. Zambon, and C. Alippi, “**Sparse graph learning from spatiotemporal time series,**” *Journal of Machine Learning Research*, vol. 24, no. 242, pp. 1–36, 2023.

# References vii

---

- [35] T. Kipf, E. Fetaya, K.-C. Wang, M. Welling, and R. Zemel, “**Neural relational inference for interacting systems,**” in *International conference on machine learning*, PMLR, 2018, pp. 2688–2697.
- [36] C. Gray, L. Mitchell, and M. Roughan, “**Bayesian inference of network structure from information cascades,**” *IEEE Transactions on Signal and Information Processing over Networks*, vol. 6, pp. 371–381, 2020.
- [37] A. Manenti, D. Zambon, and C. Alippi, ***Learning Latent Graph Structures and their Uncertainty***, May 2024.
- [38] J. Hosking, “**Equivalent Forms of the Multivariate Portmanteau Statistic,**” *Journal of the Royal Statistical Society: Series B (Methodological)*, vol. 43, no. 2, pp. 261–262, 1981.
- [39] Z. Li, C. Lam, J. Yao, and Q. Yao, “**On Testing for High-Dimensional White Noise,**” *The Annals of Statistics*, vol. 47, no. 6, pp. 3382–3412, 2019.
- [40] A. Bose and W. Hachem, “**A Whiteness Test Based on the Spectral Measure of Large Non-Hermitian Random Matrices,**” in *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, IEEE, 2020, pp. 8768–8771.
- [41] P. A. P. Moran, “**Notes on Continuous Stochastic Phenomena,**” *Biometrika*, vol. 37, no. 1/2, pp. 17–23, 1950, ISSN: 0006-3444. DOI: 10.2307/2332142.

## References viii

---

- [42] A. D. Cliff and K. Ord, “**Spatial Autocorrelation: A Review of Existing and New Measures with Applications,**” *Economic Geography*, vol. 46, pp. 269–292, 1970, ISSN: 0013-0095. DOI: 10.2307/143144.
- [43] D. Zambon and C. Alippi, “**AZ-whiteness test: A test for signal uncorrelation on spatio-temporal graphs,**” in *Advances in Neural Information Processing Systems*, A. H. Oh, A. Agarwal, D. Belgrave, and K. Cho, Eds., 2022.
- [44] D. Zambon and C. Alippi, “**Where and how to improve graph-based spatio-temporal predictors,**” *arXiv preprint arXiv:2302.01701*, 2023.
- [45] B. Yu, H. Yin, and Z. Zhu, “**ST-Unet: A spatio-temporal U-network for graph-structured time series modeling,**” *arXiv preprint arXiv:1903.05631*, 2019.
- [46] A. Cini, D. Mandic, and C. Alippi, “**Graph-based Time Series Clustering for End-to-End Hierarchical Forecasting,**” *International Conference on Machine Learning*, 2024.
- [47] S. S. Rangapuram, M. W. Seeger, J. Gasthaus, L. Stella, Y. Wang, and T. Januschowski, “**Deep State Space Models for Time Series Forecasting,**” in *Advances in Neural Information Processing Systems*, vol. 31, Curran Associates, Inc., 2018.

# References ix

---

- [48] D. Zambon, A. Cini, L. Livi, and C. Alippi, **Graph state-space models**, Jan. 2023. doi: 10.48550/arXiv.2301.01741.
- [49] C. Alippi and D. Zambon, **Graph Kalman Filters**, Mar. 2023. doi: 10.48550/arXiv.2303.12021.
- [50] I. Buchnik, G. Sagi, N. Leinwand, Y. Loya, N. Shlezinger, and T. Routtenberg, **“Gsp-kalmanet: Tracking graph signals via neural-aided kalman filtering,”** *IEEE Transactions on Signal Processing*, 2024.
- [51] E. Chouzenoux and V. Elvira, **“Sparse graphical linear dynamical systems,”** *Journal of Machine Learning Research*, vol. 25, no. 223, pp. 1–53, 2024.
- [52] X. Yin, F. Li, Y. Shen, H. Qi, and B. Yin, **“Nodetrans: A graph transfer learning approach for traffic prediction,”** *arXiv preprint arXiv:2207.01301*, 2022.
- [53] A. Prabowo, H. Xue, W. Shao, P. Koniusz, and F. D. Salim, **“Traffic forecasting on new roads using spatial contrastive pre-training (scpt),”** *Data Mining and Knowledge Discovery*, vol. 38, no. 3, pp. 913–937, 2024.

# References x

---

- [54] W. Hu, M. Fey, M. Zitnik, *et al.*, “**Open graph benchmark: Datasets for machine learning on graphs,**” in *Advances in Neural Information Processing Systems*, H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan, and H. Lin, Eds., vol. 33, Curran Associates, Inc., 2020, pp. 22 118–22 133. [Online]. Available: [https://proceedings.neurips.cc/paper\\_files/paper/2020/file/fb60d411a5c5b72b2e7d3527cfc84fd0-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/fb60d411a5c5b72b2e7d3527cfc84fd0-Paper.pdf).
- [55] S. Huang, F. Poursafaei, J. Danovitch, *et al.*, “**Temporal graph benchmark for machine learning on temporal graphs,**” in *Advances in Neural Information Processing Systems*, A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, Eds., vol. 36, Curran Associates, Inc., 2023, pp. 2056–2073. [Online]. Available: [https://proceedings.neurips.cc/paper\\_files/paper/2023/file/066b98e63313162f6562b35962671288-Paper-Datasets\\_and\\_Benchmarks.pdf](https://proceedings.neurips.cc/paper_files/paper/2023/file/066b98e63313162f6562b35962671288-Paper-Datasets_and_Benchmarks.pdf).
- [56] J. Gastinger, S. Huang, M. Galkin, *et al.*, “**Tgb 2.0: A benchmark for learning on temporal knowledge graphs and heterogeneous graphs,**” in *Advances in Neural Information Processing Systems*, 2024.
- [57] X. Liu, Y. Xia, Y. Liang, *et al.*, “**Largest: A benchmark dataset for large-scale traffic forecasting,**” *Advances in Neural Information Processing Systems*, vol. 36, 2024.



# References xi

---

- [58] Z. Shao, F. Wang, Y. Xu, *et al.*, “**Exploring progress in multivariate time series forecasting: Comprehensive benchmarking and heterogeneity analysis,**” *IEEE Transactions on Knowledge and Data Engineering*, pp. 1–14, 2024. DOI: 10.1109/TKDE.2024.3484454.