

Learning on Graphs Conference Italy meetup · December 6, 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

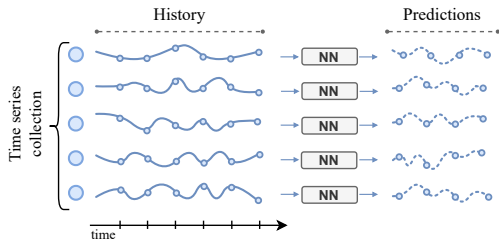


Stock markets

# 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**.



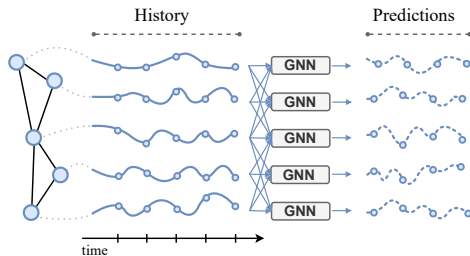
[1] Salinas *et al.*, “DeepAR: Probabilistic forecasting with autoregressive recurrent networks”, IJF 2020.

[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

---

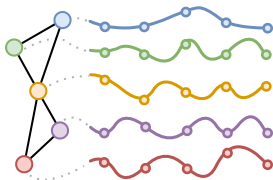
This tutorial presents advances coming from the **combination** of

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

The **objective** of this **short** tutorial is to provide:

1. a framework for **graph-based time series processing** models;
2. a discussion of selected **challenges** and **future directions**.

There is a **longer version** of this tutorial<sup>1</sup>, complemented by a **software demo** and a **paper** [3].



---

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

<sup>1</sup>Available at [gmlg.ch](http://gmlg.ch)

Part 1

# Graph-based Processing of Correlated Time series

## **Correlated time series**

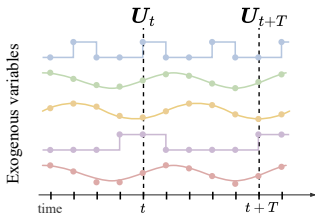
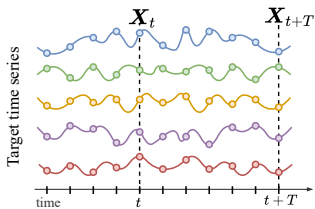
---



# 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.,  $\mathbf{X}_t \in \mathbb{R}^{N \times d_x}$ ,  $\mathbf{U}_t \in \mathbb{R}^{N \times d_u}$ .

→ 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:**

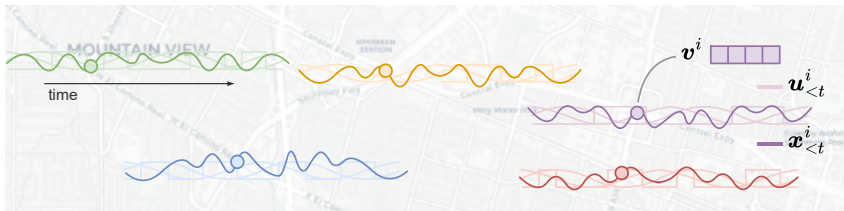
$$\mathcal{X}_t = \langle \mathbf{X}_t, \mathbf{U}_t, \mathbf{V} \rangle$$

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

! 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.



- $X_{<t}$  collects past traffic speed measurements.
- $U_t$  stores identifiers for time-of-the-day and day-of-the-week.
- $V$  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

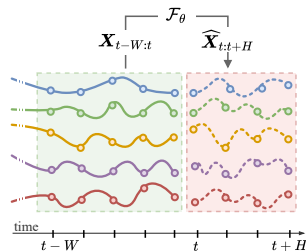
## Multi-step time-series forecasting

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.

# Global and local predictors

## Local models

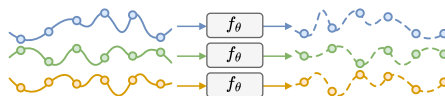


$$\hat{x}_{t+h}^i = f(x_{t-W:t}^i; \theta^i)$$

**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

---

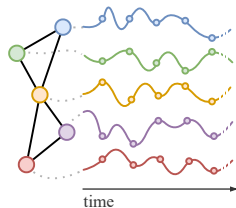


# Relational information

💡 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  $t$ ).

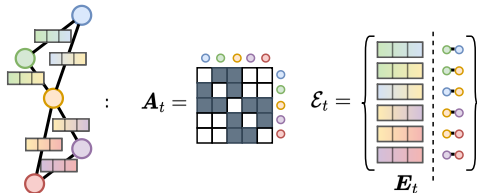


# Relational information with attributes

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**.

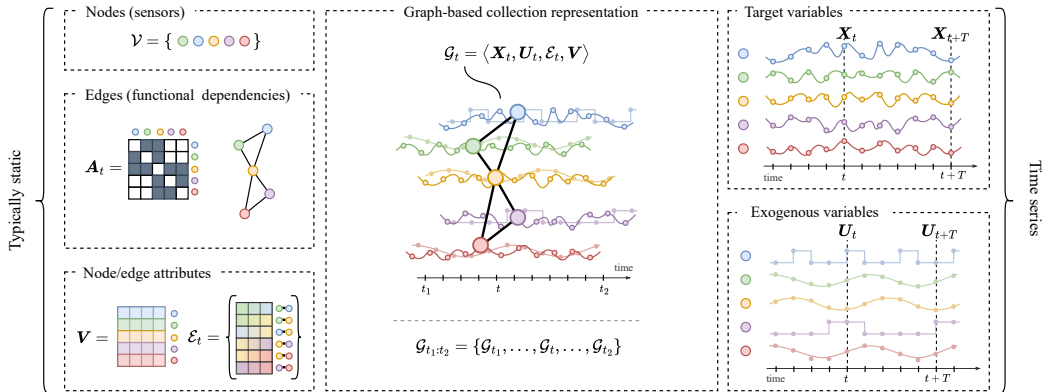
# Example: Traffic monitoring system

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



$\mathcal{G}_t \doteq \langle X_t, U_t, E_t, V \rangle$  contains the available information w.r.t. time step  $t$ .

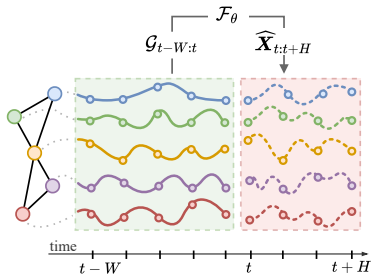
# Relational inductive biases for time series forecasting

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

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

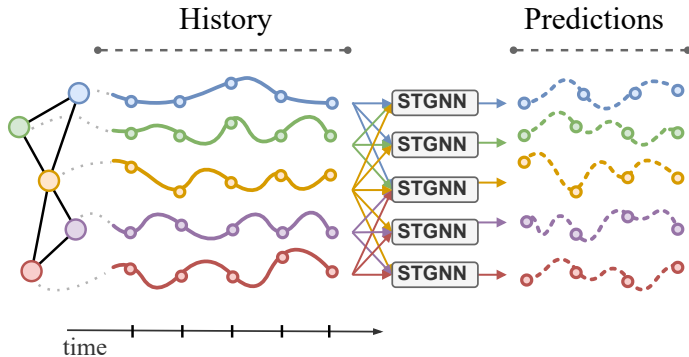
The conditioning can act as a **regularization** to localize predictions w.r.t. **each node**.

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



# Spatiotemporal graph neural networks

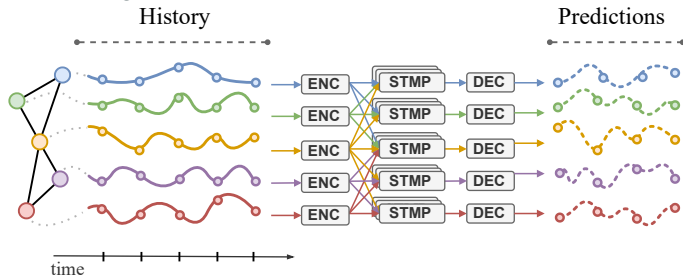
We call [spatiotemporal graph neural networks \(STGNNs\)](#) a neural network exploiting both temporal and spatial relations of the input spatiotemporal time series.



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.

# 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.

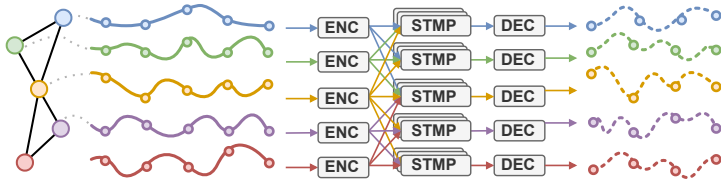
[6] Gilmer *et al.*, “Neural message passing for quantum chemistry”, ICML 2017.



# **Globality and locality in STGNNs**

---

# Globality and locality in STGNNs

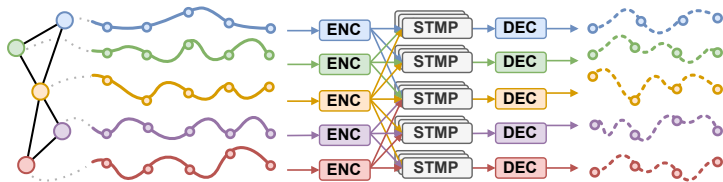


Standard STGNNs are **global** models.

- 😊 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.

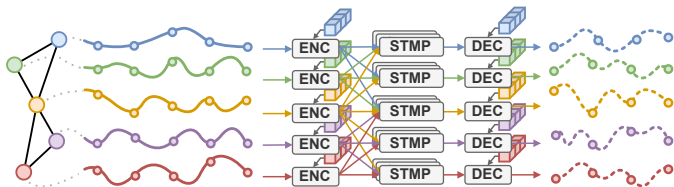
# Global-local STGNNs



💡 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**.

- 😊 Most of the model's parameters remain shared.
- 😊 Can facilitate transfer learning.
- 😞 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.

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

# What we have seen so far

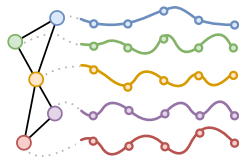
---

1. Introduced the problem of processing **correlated time series**.
2. **Graph representations** allows for modeling dependencies.
3. Discussed the **forecasting problem** and associated **predictors**.
4. Saw recipes for building (global/local) **STGNNs**.

In the following, we will look into

- dealing with partial observations;
- latent graph learning;
- a selection of future directions.

Checkout the full tutorial for more on: [computational scalability](#), [model quality assessment](#), [software libraries](#), . . .



Part 2

# Challenges

## **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.

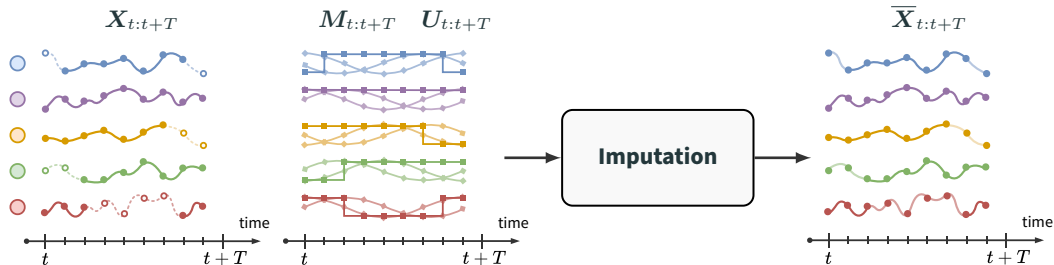
→ We need a way to **impute**, i.e., *reconstruct*, missing data.



# 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:

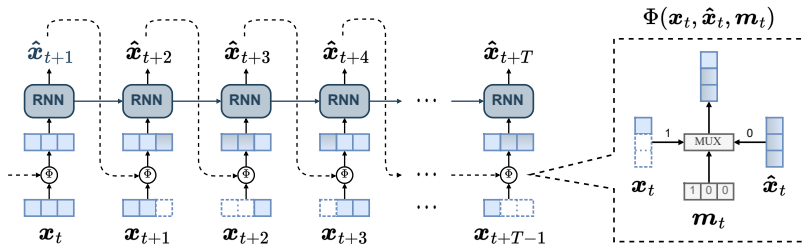
$$\hat{\theta} = \arg \min_{\theta} \sum_{t=1}^T \sum_{i=1}^N \frac{\|\overline{\mathbf{m}}_t^i \odot \ell(\overline{\mathbf{x}}_t^i, \mathbf{x}_t^i)\|_1}{\|\overline{\mathbf{m}}_t^i\|_1}.$$

 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).

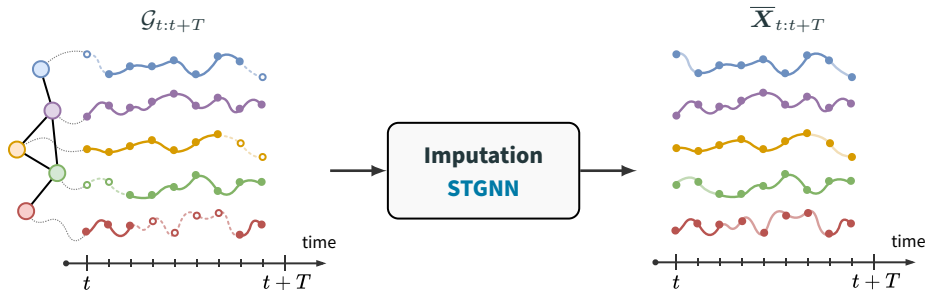


- 😊 Effective in exploiting past (and future, with bidirectional models) **node** observations.
- 😞 Struggle in capturing **nonlinear space-time dependencies**.

# 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 [8].

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.

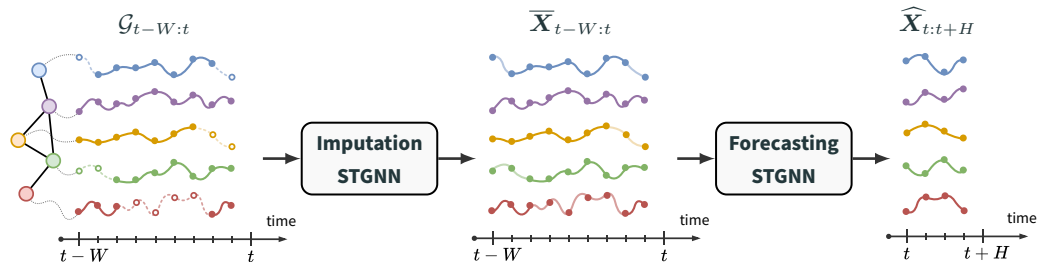
$$p(\mathbf{x}_t^i | \{\mathbf{x}_t^j \odot \mathbf{m}_t^j\}^{j \neq i})$$

Enabled by message passing.

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

# Imputation *before* forecasting

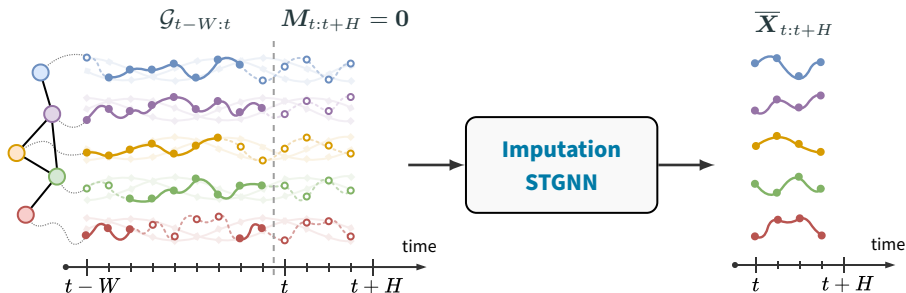
TSI is often used as a **preprocessing step** for a downstream task, e.g., forecasting.



- ☹️ 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.



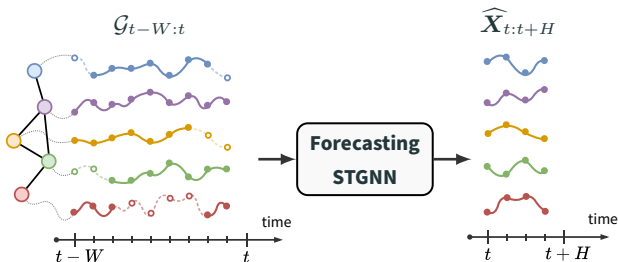
# Forecasting from partial observations

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.**



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

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

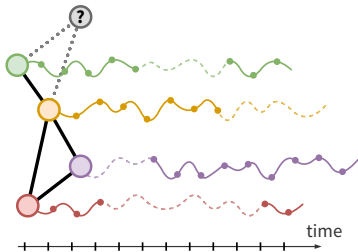
[11] Marisca *et al.*, “Graph-based Forecasting with Missing Data through Spatiotemporal Downsampling”, ICML 2024.

# 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.



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

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

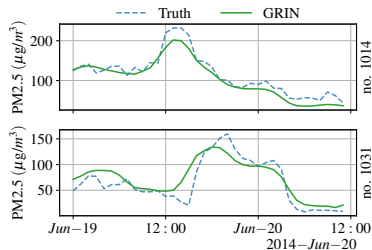
# Graph imputation for virtual sensing

💡 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 [8])



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

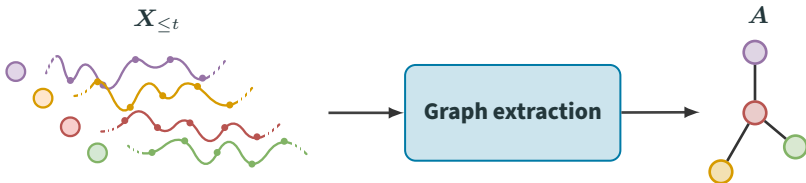
[14] Marisca *et al.*, “Learning to Reconstruct Missing Data from Spatiotemporal Graphs with Sparse Observations”, NeurIPS 2022.

# 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



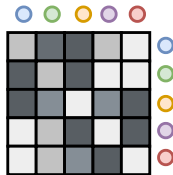
- When possible, the learned graph should be **sparse**.
- It can be interpreted as **regularizing a spatial attention** operator.
- This task is found under different names:  
**graph structure learning**, latent graph learning, graph inference...

# 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.

---

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

[16] 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

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

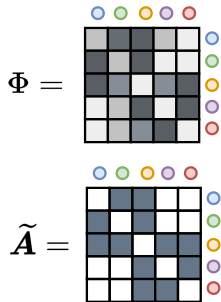
Edge scores  $\Phi$

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

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

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

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

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**.

## 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.,

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

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

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

# Learning graph distributions

---

Training typically involves optimizing terms similar to

$$\mathcal{L}(\theta, \Phi) = \mathbb{E}_{\mathbf{A} \sim p_{\Phi}} [L_{\theta}(\mathbf{A})]$$

which average a cost  $L_{\theta}$  over all graphs according to  $p_{\Phi}$ .

For example,

$$\mathcal{L}(\theta, \Phi) =$$

$$\text{(MSE)} \quad = \mathbb{E}_{\mathbf{A} \sim p_{\Phi}} [\|\mathcal{F}_{\theta}(\mathbf{X}_{t-W:t}, \mathbf{A}) - \mathbf{X}_{t:t+H}\|^2].$$

$$\text{(CRPS)} \quad = \mathbb{E}_{\mathbf{A} \sim p_{\Phi}} [\|\mathcal{F}_{\theta}(\mathbf{X}_{t-W:t}, \mathbf{A}) - \mathbf{X}_{t:t+H}\|] - \frac{1}{2} \mathbb{E}_{\mathbf{A}, \mathbf{A}' \sim p_{\Phi}} [\|\mathcal{F}_{\theta}(\mathbf{X}_{t-W:t}, \mathbf{A}) - \mathcal{F}_{\theta}(\mathbf{X}_{t-W:t}, \mathbf{A}')\|].$$

(The expected value  $\mathbb{E}_{\mathbf{X}, \mathbf{Y}}$  over the input and output data distribution is omitted for brevity.)

# Gradient-based optimization and Monte Carlo sampling

Gradient-based optimization requires the computation of  $\nabla_{\theta}$  and  $\nabla_{\Phi}$  of  $\mathcal{L}(\theta, \Phi) = \mathbb{E}_{\mathbf{A} \sim p_{\Phi}} [L_{\theta}(\mathbf{A})]$ .

😊 Gradient  $\nabla_{\theta} \mathcal{L}(\theta, \Phi)$  can be estimated via Monte Carlo (MC) with standard tools

$$\nabla_{\theta} \mathcal{L}(\theta, \Phi) \stackrel{MC}{\approx} \nabla_{\theta} \frac{1}{M} \sum_m L_{\theta}(\mathbf{A}^m) = \frac{1}{M} \sum_m \nabla_{\theta} L_{\theta}(\mathbf{A}^m)$$

with  $\{\mathbf{A}^m\}_{m=1}^M$  being a set of i.i.d.  $M$  samples from  $p_{\Phi}$ .

☹️ Estimating gradient  $\nabla_{\Phi} \mathcal{L}(\theta, \Phi)$  via MC is less straightforward:

$$\nabla_{\Phi} \mathcal{L}(\theta, \Phi) = \nabla_{\Phi} \mathbb{E}_{\mathbf{A} \sim p_{\Phi}} [L_{\theta}(\mathbf{A})]$$

☹️ Expanding the gradient leads to

$$\nabla_{\Phi} \mathcal{L}(\theta, \Phi) = \int L_{\theta}(\mathbf{A}) \nabla_{\Phi} p_{\Phi}(\mathbf{A}) d\mathbf{A}.$$

- not in the form of an expected value,
- analytical computation is often unfeasible.

## Reparametrization trick

---

💡 One approach is to **reparametrize**  $\mathbf{A} \sim p_{\Phi}(\mathbf{A})$  as:  $\mathbf{A} = g(\Phi, \varepsilon)$ ,  $\varepsilon \sim p(\varepsilon)$

→ for instance,  $a \sim \mathcal{N}(\mu, \sigma)$  can be written as  $a = \mu + \varepsilon\sigma$ , with  $\varepsilon \sim \mathcal{N}(0, 1)$ .

Above rewriting decouples parameters  $\Phi$  from the random component  $\varepsilon$ :

$$\nabla_{\Phi} \mathbb{E}_{\mathbf{A} \sim p_{\Phi}} [L_{\theta}(\mathbf{A})] = \mathbb{E}_{\varepsilon} [\nabla_{\Phi} L(g(\Phi, \varepsilon))].$$

If  $\mathbf{A} \in \{0, 1\}$ , gradient  $\nabla_{\Phi} g_{\Phi}(A) = 0$  almost everywhere and undefined otherwise.

→ Continuous relaxation is used, e.g., Concrete distribution.

😊 Relatively **easy** to implement,

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

---

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

[21] Elinas *et al.*, “Variational inference for graph convolutional networks in the absence of graph data and adversarial settings”, NeurIPS 2020.

# Score-function gradient estimator

💡 Score-function gradient estimators rely on the relation

$$\nabla_{\Phi} \mathbb{E}_{p_{\Phi}} [L_{\theta}(\mathbf{A})] = \mathbb{E}_{p_{\Phi}} [L_{\theta}(\mathbf{A}) \nabla_{\Phi} \log p_{\Phi}(\mathbf{A})]$$

In our forecasting settings, it reads

$$\nabla_{\Phi} \mathcal{L}(\theta, \Phi) \stackrel{MC}{\approx} \frac{1}{M} \sum_{m=1}^M \ell(\mathcal{F}_{\theta}(\mathbf{X}_{t-W:t}, \mathbf{A}), \mathbf{X}_{t:t+H}) \nabla_{\Phi} \log p_{\Phi}(\mathbf{A}_m).$$

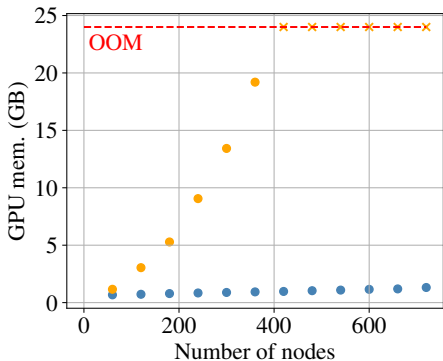
- ☹️ suffer from **high variance** (use variance reduction techniques),
- 😊 allow to **keep computations sparse** through the model.
  - do not rely on continuous relaxation of **discrete random variables**;
  - allow for **sparse message passing** in  $\mathcal{F}(\mathbf{X}_{t-W:t}, \mathbf{A})$  by relying on sparse matrices  $\mathbf{A}$ .

---

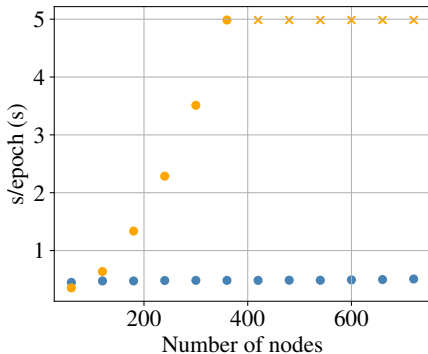
[19] Cini *et al.*, “Sparse graph learning from spatiotemporal time series”, JMLR 2023.



# Computational efficiency



● Score-function



● Reparametrization trick

# Uncertainty quantification

---

While probabilistic models have been used to enable the **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.

---

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

# Learning guarantees for latent graph calibration

😊 Under appropriate assumptions, we can achieve:  $p_{\Phi}(\widehat{\mathbf{X}}_{t:t+H} | \mathbf{X}_{t-W:t}) = p(\mathbf{X}_{t:t+H} | \mathbf{X}_{<t})$ .

→ This means that the model **output** is **calibrated**.

😞 Calibration of  $\mathbf{A}$  is not implied from that of the model output.

→ **Conditions** on the function  $\mathbf{A} \mapsto \widehat{\mathbf{X}}_{t:t+H} = \mathcal{F}_{\Phi}(\mathbf{X}_{t-W:t}, \mathbf{A})$  are **requested**.

😊 For **graphs** and graph neural networks, these conditions appear **easier to meet!**

---

[23] Gneiting *et al.*, “Probabilistic forecasting”, *Annu. Rev. Stat. Appl.* 2014.

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

Part 3

# Future Directions

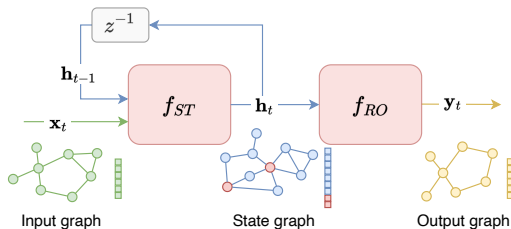
# **Graph State-Space Models**

---

# 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.



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

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

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

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

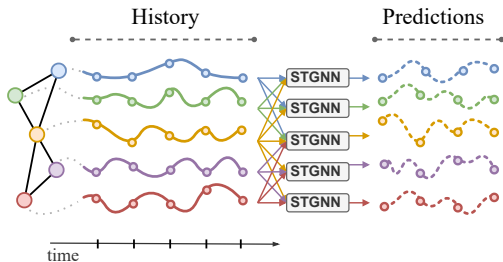
[29] Chouzenoux *et al.*, “Sparse graphical linear dynamical systems”, JMLR 2024.

# **Hierarchical processing**

---

# What we achieved so far

- 😊 Pairwise dependencies are embedded into the processing.
- 😊 Predictions are localized w.r.t. a node and its neighbors.
- 😞 Operate at a fixed spatiotemporal scale.
- 😞 Higher-order dependencies are not explicitly modeled.



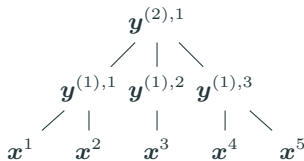


# Hierarchical forecasting

- Hierarchical forecasting is about making predictions at multiple resolutions.
- Coherency constraints provide a regularization mechanism.
- Predictions are coherent iff:

$$Q\hat{Y}_t = [I \mid -C] \hat{Y}_t = \mathbf{0},$$

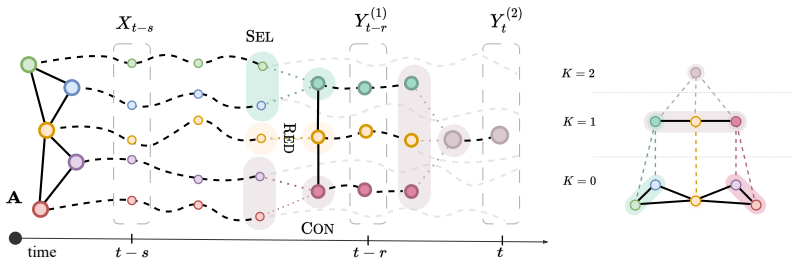
where  $\mathbf{Y}_t$  contains stacked predictions for each level.



$$C = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 1 \end{bmatrix}$$

[30] Hyndman *et al.*, “Optimal combination forecasts for hierarchical time series”, Elsevier CSDA 2011.

# Hierarchical Graph Predictor (HiGP)



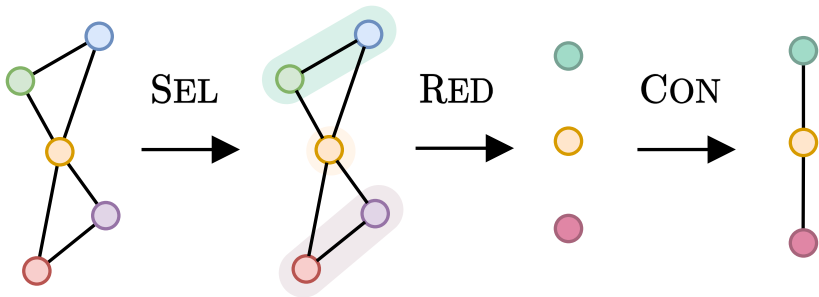
💡 Combine **hierarchical** and **graph-based** forecasting.

We introduced a framework unifying the two.

- 😊 Operates at **different spatial resolutions** exploring a pyramidal graph structure.
- 😊 Exploits **higher-order dependencies** by operating on **aggregated** time series.
- 😊 Hierarchical time series **clusters** are learned **end-to-end**.

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

# Select, Reduce, Connect (SRC)



- **Select** maps nodes into supernodes, i.e., node clusters.

- **Reduce** specifies how observations should be aggregated.

- **Connect** specifies how to rewire the graph after pooling.

[32] Grattarola *et al.*, “Understanding Pooling in Graph Neural Networks”, IEEE TNLS 2022.

## A possible implementation

---

By exploiting the SRC framework to define the proper operators, [hierarchical architectural biases](#) can be embedded into a [time-then-space STGNN](#) architecture as

$$\begin{aligned}
 \mathbf{h}_t^{(k),i,0} &= \text{SEQENC}^{(k)} \left( \mathbf{y}_{t-W:t}^{(k),i}, \mathbf{u}_{t-W:t}^{(k),i}, \mathbf{v}^{(k),i} \right), & \text{Temporal enc.} \\
 \mathbf{Z}_t^{(k),l} &= \text{MP}_l^{(k)} \left( \mathbf{H}_t^{(k),l}, \mathbf{A}^{(k)} \right), & \text{Intra-level prop.} \\
 \mathbf{H}_t^{(k),l+1} &= \text{MLP}_l^{(k)} \left( \mathbf{Z}_t^{(k),l}, \underbrace{\mathbf{S}^{(k)T} \mathbf{Z}_t^{(k-1),l}}_{\text{RED}^{(k)}}, \underbrace{\mathbf{S}^{(k)} \mathbf{Z}_t^{(k+1),l}}_{\text{LIFT}^{(k)}} \right). & \text{Inter-level prop.}
 \end{aligned}$$

Representations can then be mapped to prediction using a **readout**.

# Making coherent hierarchical forecasts

- **Learning time series clusters end-to-end**

We learn **probabilistic cluster assignments** and use a **MinCut** [33] regularize.

$$\mathbf{S}^{(k)} \sim P(\mathbf{S}_{ij}^{(k)} = 1) = \frac{e^{\phi_{ij}^{(k)}/\tau}}{\sum_j e^{\phi_{ij}^{(k)}/\tau}}, \quad \Phi^{(k)} = \mathcal{F}_\psi \left( \mathbf{Y}_{t-W:t}^{(k-1)}, \mathbf{A}^{(k-1)}, \mathbf{V}^{(k-1)} \right).$$

Forecasting the resulting aggregates provides a **self-supervised** learning mechanism.

- **Forecast reconciliation**

A differentiable reconciliation step ensures **coherent forecasts** by recombining predictions as

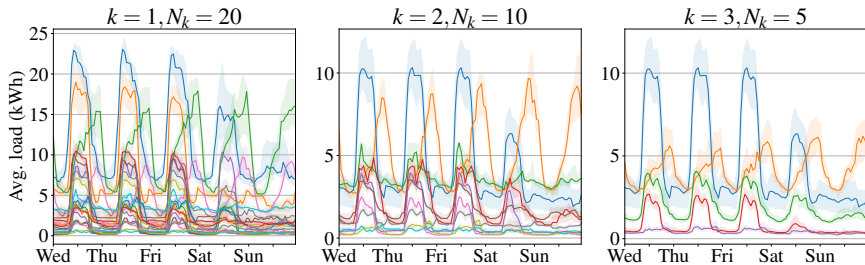
$$\mathbf{P} \doteq \mathbf{I} - \mathbf{Q}^T \left( \mathbf{Q}\mathbf{Q}^T \right)^{-1} \mathbf{Q}, \quad \bar{\mathbf{Y}}_t = \mathbf{P}\hat{\mathbf{Y}}_t.$$

☹️ Computing the inverse has a **cubic cost**, a soft regularization can be used alternatively.

[33] Bianchi *et al.*, “Spectral clustering with graph neural networks for graph pooling”, ICML 2020.

[34] Rangapuram *et al.*, “End-to-end learning of coherent probabilistic forecasts for hierarchical time series”, ICML 2021.

# Example of learned clusters



Learned hierarchical clusters from CER (energy consumption profiles).

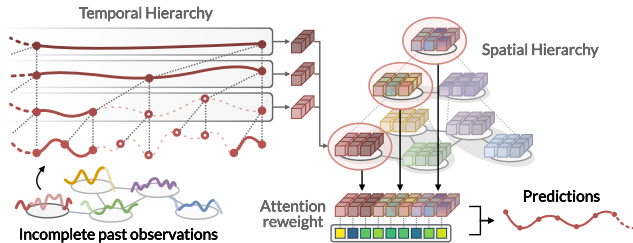
[31] Cini *et al.*, “Graph-based Time Series Clustering for End-to-End Hierarchical Forecasting”, ICML 2024.

# Multi-scale spatiotemporal representations

Similar hierarchical processing can be jointly performed also over the **temporal** dimension.

This gives a **hierarchy of multi-scale representations**, each accounting for a specific **space-time resolution**.

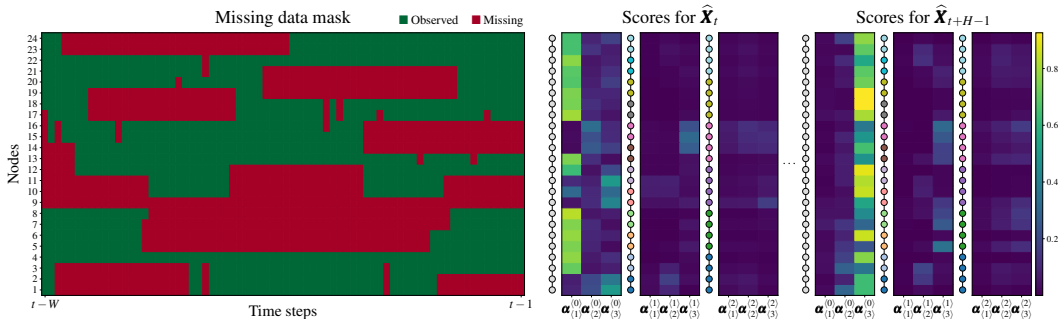
- 😊 Different scales might capture **different dynamics**.
- 😊 Helps with **noisy and missing data**.



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

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

# Downsampling with missing data



- The model focuses on the **fine-grained temporal scale** – if the most recent data are **not missing**.
- When data are missing at a given node, **higher levels** in the **spatial hierarchy** are given more weight.
- **Slower dynamics** become more relevant when **long-range** forecasting.

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

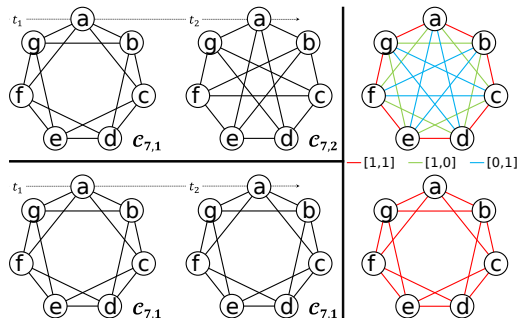


# Theoretical properties

---

# Theoretical properties of STGNNs

- High interest in studying the **expressive power** of GNNs in *static* graphs [36].
- Recent work extended the focus to dynamic settings, e.g., **temporal graphs**.
- **The important question**  
What's the impact of **different spatiotemporal message-passing operators** on the properties of the resulting STGNN?



From [37].

[37] Gao *et al.*, “On the Equivalence Between Temporal and Static Equivariant Graph Representations”, ICML 2022.

[38] Gravina *et al.*, “Long Range Propagation on Continuous-Time Dynamic Graphs”, ICML 2024.

[39] Beddar-Wiesing *et al.*, “Weisfeiler–Lehman goes dynamic: An analysis of the expressive power of graph neural networks for attributed and dynamic graphs”, Neural Networks 2024.

[40] Wałęga *et al.*, “Expressive Power of Temporal Message Passing”, Preprint 2024.

# Conclusions

---

## Some Takeaways

---

Deep Learning  
for **time series**

+

Deep Learning  
on **graphs**

- ⚡ Relational inductive **biases** allow for exploiting dependencies among the time series...
- 😊 ... and **effectively processing** spatiotemporal data,
- 😊 while **sharing** most of the model **parameters**.
- 💡 **Global-local models** are a good starting point.

**Resources.** 📄 Tutorial paper [3] • 🗄️ Open-source library [41]

---

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

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



**Andrea Cini**



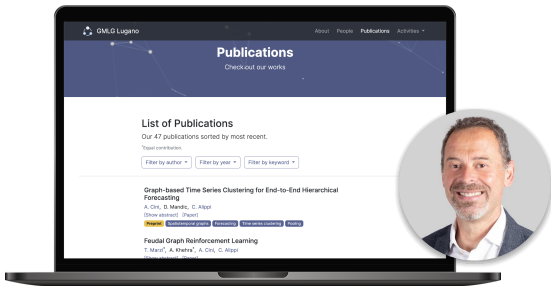
**Ivan Marisca**



**Daniele Zambon**

**Graph Machine Learning Group**  
**gmlg.ch**

Group leader: Prof. Cesare Alippi



# 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] A. Cini, I. Marisca, D. Zambon, and C. Alippi, “**Taming local effects in graph-based spatiotemporal forecasting,**” *arXiv preprint arXiv:2302.04071*, 2023.
- [8] 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>.
- [9] X. Zhang, M. Zeman, T. Tsiligkaridis, and M. Zitnik, “**Graph-guided network for irregularly sampled multivariate time series,**”, 2022.
- [10] 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.
- [11] 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 iii

---

- [12] 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.
- [13] 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>.
- [14] 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.
- [15] 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.
- [16] 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.
- [17] 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 iv

---

- [18] 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.
- [19] 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.
- [20] 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.
- [21] P. Elinas, E. V. Bonilla, and L. Tiao, “**Variational inference for graph convolutional networks in the absence of graph data and adversarial settings,**” *Advances in neural information processing systems*, vol. 33, pp. 18 648–18 660, 2020.
- [22] 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.
- [23] T. Gneiting and M. Katzfuss, “**Probabilistic forecasting,**” *Annual Review of Statistics and Its Application*, vol. 1, no. 1, pp. 125–151, 2014.

# References v

---

- [24] A. Manenti, D. Zambon, and C. Alippi, ***Learning Latent Graph Structures and their Uncertainty***, May 2024.
- [25] 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.
- [26] D. Zambon, A. Cini, L. Livi, and C. Alippi, ***Graph state-space models***, Jan. 2023. doi: 10.48550/arXiv.2301.01741.
- [27] C. Alippi and D. Zambon, ***Graph Kalman Filters***, Mar. 2023. doi: 10.48550/arXiv.2303.12021.
- [28] I. Buchnik, G. Sagi, N. Leinwand, Y. Loya, N. Shlezinger, and T. Rauttenberg, **“Gsp-kalmanet: Tracking graph signals via neural-aided kalman filtering,”** *IEEE Transactions on Signal Processing*, 2024.
- [29] E. Chouzenoux and V. Elvira, **“Sparse graphical linear dynamical systems,”** *Journal of Machine Learning Research*, vol. 25, no. 223, pp. 1–53, 2024.
- [30] R. J. Hyndman, R. A. Ahmed, G. Athanasopoulos, and H. L. Shang, **“Optimal combination forecasts for hierarchical time series,”** *Computational statistics & data analysis*, vol. 55, no. 9, pp. 2579–2589, 2011.

# References vi

---

- [31] A. Cini, D. Mandic, and C. Alippi, “**Graph-based Time Series Clustering for End-to-End Hierarchical Forecasting,**” *International Conference on Machine Learning*, 2024.
- [32] D. Grattarola, D. Zambon, F. Bianchi, and C. Alippi, “**Understanding Pooling in Graph Neural Networks,**” *IEEE Transactions on Neural Networks and Learning Systems*, pp. 1–11, 2022. DOI: 10.1109/TNNLS.2022.3190922.
- [33] F. M. Bianchi, D. Grattarola, and C. Alippi, “**Spectral clustering with graph neural networks for graph pooling,**” in *International conference on machine learning*, PMLR, 2020, pp. 874–883.
- [34] S. S. Rangapuram, L. D. Werner, K. Benidis, P. Mercado, J. Gasthaus, and T. Januschowski, “**End-to-end learning of coherent probabilistic forecasts for hierarchical time series,**” in *International Conference on Machine Learning*, PMLR, 2021, pp. 8832–8843.
- [35] 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.
- [36] K. Xu, W. Hu, J. Leskovec, and S. Jegelka, “**How powerful are graph neural networks?**” In *International Conference on Learning Representations*, 2019.

# References vii

---

- [37] J. Gao and B. Ribeiro, “**On the equivalence between temporal and static equivariant graph representations,**” in *International Conference on Machine Learning*, PMLR, 2022, pp. 7052–7076.
- [38] A. Gravina, G. Lovisotto, C. Gallicchio, D. Bacciu, and C. Grohnfeldt, “**Long range propagation on continuous-time dynamic graphs,**” in *Forty-first International Conference on Machine Learning*, 2024. [Online]. Available: <https://openreview.net/forum?id=gVg8V9isul>.
- [39] S. Beddar-Wiesing, G. A. D’Inverno, C. Graziani, *et al.*, “**Weisfeiler–lehman goes dynamic: An analysis of the expressive power of graph neural networks for attributed and dynamic graphs,**” *Neural Networks*, vol. 173, p. 106 213, 2024.
- [40] P. A. Wałęga and M. Rawson, “**Expressive power of temporal message passing,**” *arXiv preprint arXiv:2408.09918*, 2024.
- [41] A. Cini and I. Marisca, ***Torch Spatiotemporal***, Mar. 2022. [Online]. Available: <https://github.com/TorchSpatiotemporal/tsl>.