

# Graph Deep Learning for Spatiotemporal Time Series

Forecasting, Reconstruction and Analysis

---

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

ECML/PKDD, Turin · September 22, 2023

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

The Swiss AI Lab IDSIA

Università della Svizzera italiana





Traffic monitoring



Smart cities



Energy analytics

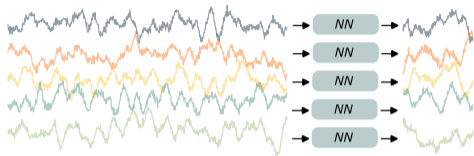


Physics



Stock markets

# Deep learning for time series forecasting



The standard deep learning approach to time series forecasting consists in training a **single neural network** on a collection of **time series**.

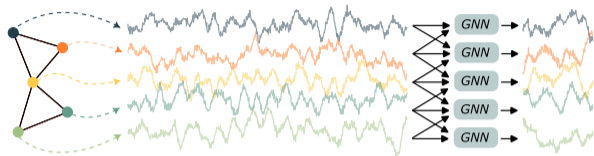
- Each time series is treated **independently** from the others.
- A single set of **shared** learnable **parameters** is used to predict each time series.
- Resulting models are **effective** and **efficient**.

**⚠ Dependencies** across time series are often **discarded**.

---

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

# Relational inductive biases



One way out is to embed such **relational structure** as an **architectural bias** into the processing.

**Graph neural networks** provide appropriate neural operators.

- **Message-passing** blocks allow for **localizing** the **predictions**  
→ conditioning on observations at related time series (neighboring nodes).
- **Parameters** are **shared** and the model can operate on arbitrary sets of time series.

---

[2] D. Bacciu *et al.*, “A gentle introduction to deep learning for graphs”, NN 2020.

[3] M. M. Bronstein *et al.*, “Geometric deep learning: Grids, groups, graphs, geodesics, and gauges” 2021.

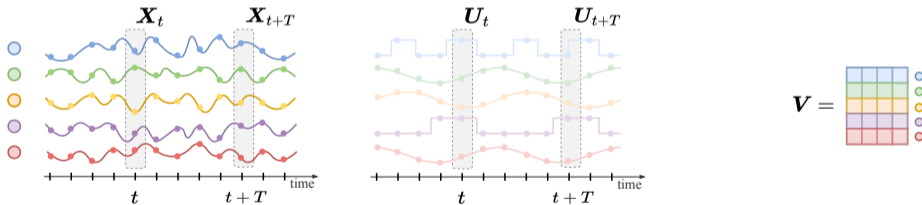
# **Spatiotemporal time series**

---

# Collections of time series

We consider a set of  $N$  **correlated time series**, where each  $i$ -th time series is associated with:

- an **observation vector**  $\mathbf{x}_t^i \in \mathbb{R}^{d_x}$  at each time step  $t$ ;
- a vector of **exogenous variable**  $\mathbf{u}_t^i \in \mathbb{R}^{d_u}$  at each time step  $t$ ;
- a vector of **static (time-independent) attributes**  $\mathbf{v}^i \in \mathbb{R}^{d_v}$ .



Capital letters denote the stacked representations encompassing the  $N$  time series in the collection, e.g.,  $\mathbf{X}_t \in \mathbb{R}^{N \times d_x}$ ,  $\mathbf{U}_t \in \mathbb{R}^{N \times d_u}$ .

[4] A. Cini *et al.*, “Graph Deep Learning for Time Series Forecasting: A Comprehensive Methodological Framework” 2023.

# Correlated time series

We assume a **time-invariant** stochastic process

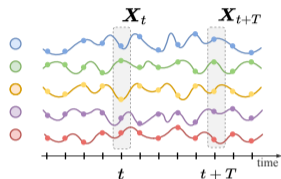
$$\mathbf{x}_t^i \sim p^i(\mathbf{x}_t^i | \mathbf{X}_{<t}, \mathbf{U}_{\leq t}, \mathbf{V})$$

generating the data  $\mathbf{x}_t^i$  for all  $i = 1 \dots N$  and  $t \in \mathbb{N}$ .

Note that the time series:

- can be generated by **different processes**,
- can **depend on each other**,
- are assumed  
homogenous, synchronous, regularly sampled.

→ These assumptions can be **relaxed**, as we will discuss in the 2nd part.



Notation:

$$\mathbf{X}_{t:t+T} = [\mathbf{X}_t, \dots, \mathbf{X}_{t+T-1}]$$

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

# Relational information

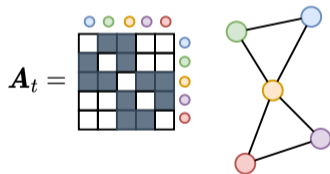
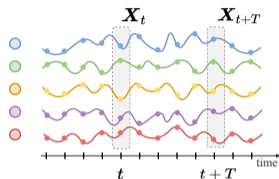
We assume the existence of **functional dependencies** between the time series.

- e.g., forecasts for one time series can be improved by accounting for the past values of other time series.

We 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  $t$ ).

- We call **spatial** the dimension spanning the time series collection.



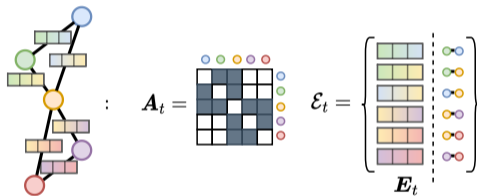


# Relational information with attributes

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

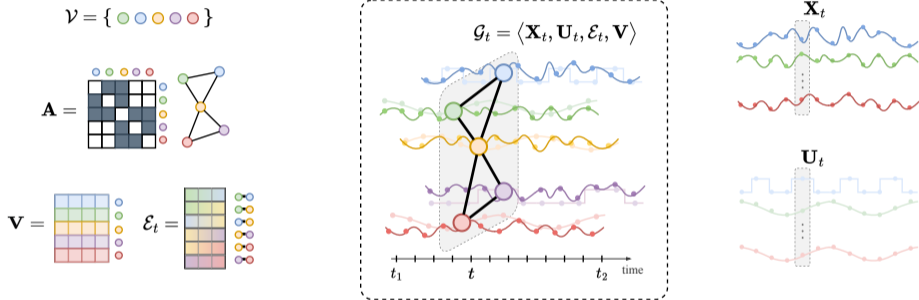
The **set of attributed edges** encoding all the available relational information is denoted by

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



→ For many applications,  $\mathbf{A}_t$  **changes slowly** over time and can be considered as **constant** within a short window of observations.

# Spatiotemporal time series



We use the terms **node** and **sensor** to indicate the  $N$  entities generating the time series.

→ We refer to the node set together with the relational information as **sensor network**.

The tuple  $\mathcal{G}_t \doteq \langle \mathbf{X}_t, \mathbf{U}_t, \mathcal{E}_t, \mathbf{V} \rangle$  contain all the available information associated with time step  $t$ .

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

# Time series forecasting

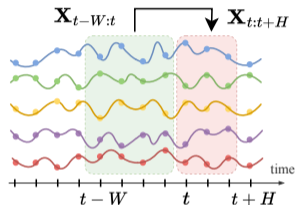
## Multi-step time-series forecasting

Given a window of  $W \geq 1$  **past** observations

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

predict  $H \geq 1$  **future** observations

$$\mathbf{x}_{t+h}^i, \quad i = 1 \dots N, h = 1 \dots H.$$



In particular, we are interested in learning a **parametric model**  $p_{\theta}$  approximating the unknown data distribution  $p$

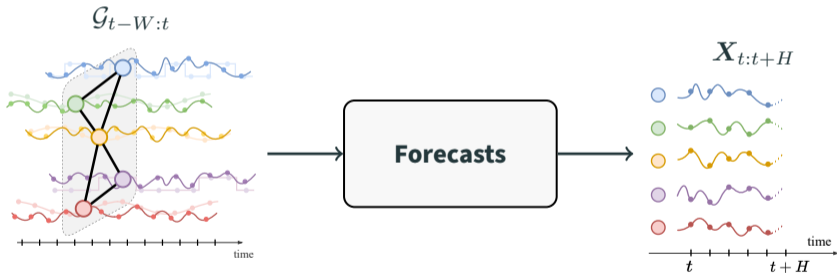
$$p_{\theta}(\mathbf{x}_{t+h}^i | \mathbf{X}_{t-W:t}, \mathbf{U}_{t-W:t+h}, \mathbf{V}) \approx p^i(\mathbf{x}_{t+h}^i | \mathbf{X}_{<t}, \mathbf{U}_{\leq t+h}, \mathbf{V}).$$

- $\theta$  is the model parameter vector.

# Time series forecasting + relational inductive biases

Condition the model on the relational information  $\mathcal{E}_{t-W:t}$

$$p_{\theta}(\mathbf{x}_{t+h}^i | \mathcal{G}_{t-W:t}, \mathbf{U}_{t-W:t+h}, \mathbf{V})$$



⚡ The conditioning on the sequence of attributed graphs acts as a **regularization** to localize predictions w.r.t. the **neighborhood of each node**.

## Point forecasts

---

For simplicity, we focus here on **point forecasts**, rather than the modeling of full data distributions  $p$ , and consider predictive model

$$\hat{\mathbf{x}}_{t+h}^i = \mathcal{F}(\mathcal{G}_{t-W:t}, \mathbf{U}_{t:t+h}; \boldsymbol{\theta})$$

where  $\hat{\mathbf{x}}_{t+h}^i$  approximates, e.g.,  $\mathbb{E}_p[\mathbf{x}_{t+h}^i]$ .

Parameters  $\boldsymbol{\theta}$  can be learned by **minimizing a cost function**  $\ell(\cdot, \cdot)$  (e.g., MSE) on a training set

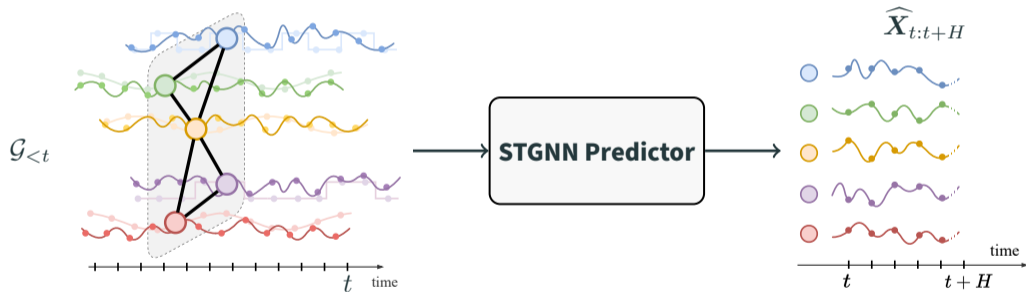
$$\begin{aligned}\hat{\boldsymbol{\theta}} &= \arg \min_{\boldsymbol{\theta}} \frac{1}{NT} \sum_{t=1}^T \ell(\hat{\mathbf{X}}_{t:t+H}, \mathbf{X}_{t:t+H}) \\ &= \arg \min_{\boldsymbol{\theta}} \frac{1}{NT} \sum_{t=1}^T \left\| \mathbf{X}_{t:t+H} - \hat{\mathbf{X}}_{t:t+H} \right\|_2^2.\end{aligned}$$

# **Spatiotemporal Graph Neural Networks**

---

# Spatiotemporal Graph Neural Networks

We call **Spatiotemporal Graph Neural Network (STGNN)** a neural network exploiting both temporal and spatial relations of the input spatiotemporal time series.



We focus on models based on **message passing**.



## Message-passing neural networks

---

To process the spatial dimension, we rely on the **message-passing (MP)** framework

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

Where:

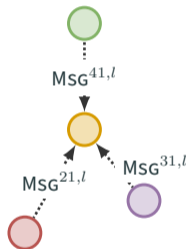
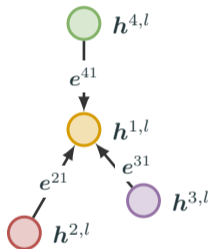
- $\text{MSG}^l(\cdot)$  is the **message function**, e.g., implemented by an MLP.
- $\text{AGGR}\{\cdot\}$  is the permutation invariant **aggregation function**.
- $\text{UP}^l(\cdot)$  is the **update function**, e.g., implemented by an MLP.

Aggregation is performed over  $\mathcal{N}(i)$ , i.e., the set of neighbors of node  $i$ .

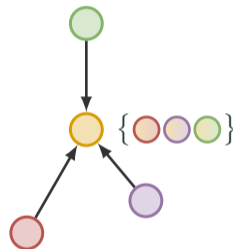
---

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

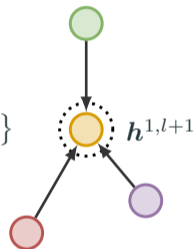
# Message passing in action



Message



Aggregate



Update

## Spatiotemporal message passing

---

Starting from the MP framework, we can define a general scheme for **spatiotemporal message-passing (STMP)** networks:

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

Rather than vectors, STMP blocks process **sequences**.

→ STMP blocks must be implemented with operators that work on sequences!

We will look at **different implementations** of STMP blocks in the following.

---

[4] A. Cini *et al.*, “Graph Deep Learning for Time Series Forecasting: A Comprehensive Methodological Framework” 2023.

## A general recipe

---

We consider STGNNs can be expressed as a sequence of three operations:

$$\mathbf{h}_{t-1}^{i,0} = \text{ENCODER}(\mathbf{x}_{t-1}^i, \mathbf{u}_{t-1}^i, \mathbf{v}^i), \quad (2)$$

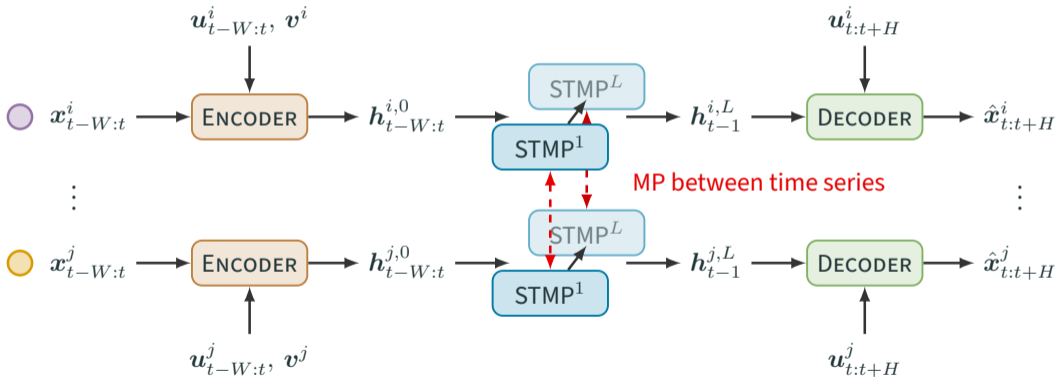
$$\mathbf{H}_{t-1}^{l+1} = \text{STMP}^l(\mathbf{H}_{\leq t-1}^l, \mathcal{E}_{\leq t-1}), \quad l = 0, \dots, L-1 \quad (3)$$

$$\hat{\mathbf{x}}_{t:t+H}^i = \text{DECODER}(\mathbf{h}_{t-1}^{i,L}, \mathbf{u}_{t:t+H}^i). \quad (4)$$

Where:

- $\text{ENCODER}(\cdot)$  is the encoding layer, e.g., implemented by an MLP.
- STMP is a stack of STMP layers.
- $\text{DECODER}(\cdot)$  is the readout layer, e.g., implemented by an MLP.

# Framework overview



## 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)**

Embed each time series in a vector, which is then propagated over the graph.

- **Space-then-Time (STT)**

Propagate nodes features at first and then process the resulting time series.



## Time-and-Space

---

In T&S models, representations at every node and time step are the results of a **joint temporal and spatial encoding**

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

Several options exist.

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

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

---

Consider a standard GRU [6] cell.

$$\mathbf{r}_t^i = \sigma (\Theta_r [\mathbf{x}_t^i || \mathbf{h}_{t-1}^i] + \mathbf{b}_r) \quad (5)$$

$$\mathbf{u}_t^i = \sigma (\Theta_u [\mathbf{x}_t^i || \mathbf{h}_{t-1}^i] + \mathbf{b}_u) \quad (6)$$

$$\mathbf{c}_t^i = \tanh (\Theta_c [\mathbf{x}_t^i || \mathbf{r}_t^i \odot \mathbf{h}_{t-1}^i] + \mathbf{b}_c) \quad (7)$$

$$\mathbf{h}_t^i = (1 - \mathbf{u}_t^i) \odot \mathbf{c}_t^i + \mathbf{u}_t^i \odot \mathbf{h}_{t-1}^i \quad (8)$$

Time series can be processed **independently** for each node or as a **single multivariate** time series.

---

[6] J. 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:

$$\mathbf{Z}_t^l = \mathbf{H}_t^{l-1} \quad (9)$$

$$\mathbf{R}_t^l = \sigma(\text{MP}_r^l([\mathbf{Z}_t^l || \mathbf{H}_{t-1}^l], \mathcal{E}_t)), \quad (10)$$

$$\mathbf{O}_t^l = \sigma(\text{MP}_o^l([\mathbf{Z}_t^l || \mathbf{H}_{t-1}^l], \mathcal{E}_t)), \quad (11)$$

$$\mathbf{C}_t^l = \tanh(\text{MP}_c^l([\mathbf{Z}_t^l || \mathbf{R}_t^l \odot \mathbf{H}_{t-1}^l], \mathcal{E}_t)), \quad (12)$$

$$\mathbf{H}_t^l = \mathbf{O}_t^l \odot \mathbf{H}_{t-1}^l + (1 - \mathbf{O}_t^l) \odot \mathbf{C}_t^l, \quad (13)$$

These T&S models are known as [graph convolutional recurrent neural networks \(GCRNNs\)](#) [7].

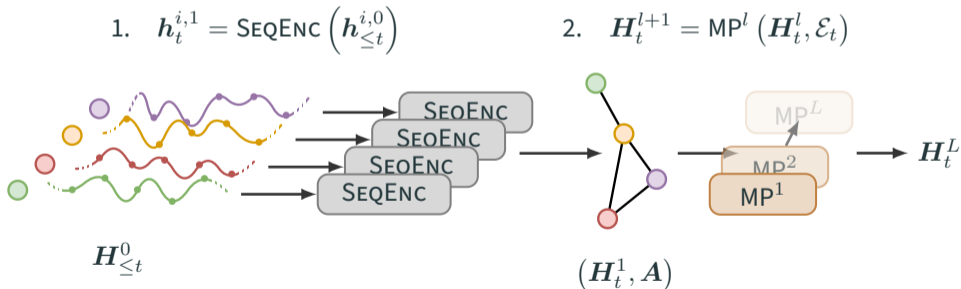
---

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

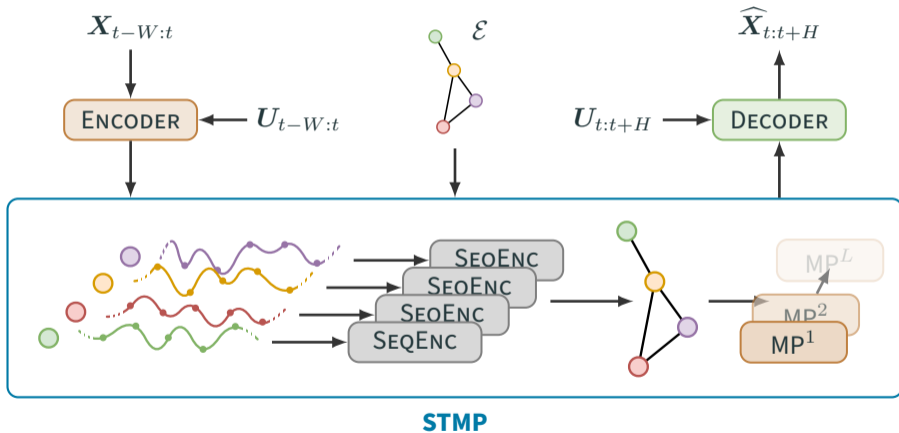
# Time-then-Space models

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.



# Full TTS model



## Pros & Cons of TTS models

---

**Pros:** 😊 Easy to implement and computationally efficient.

😊 We can reuse operators we already know.

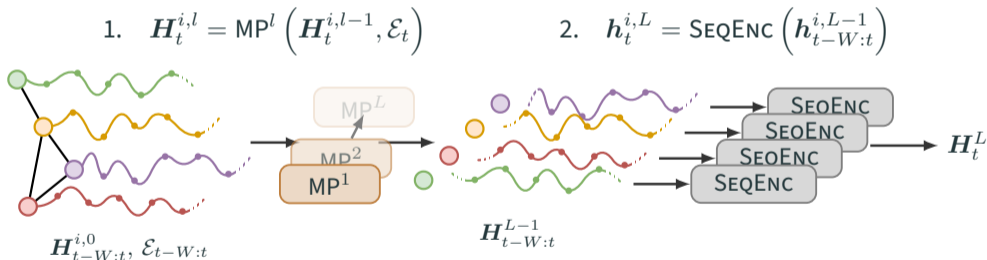
**Cons:** 😞 The 2-step encoding might introduce information bottlenecks.

😞 Accounting for changes in topology and dynamic edge attributes can be more problematic.

# Space-then-Time

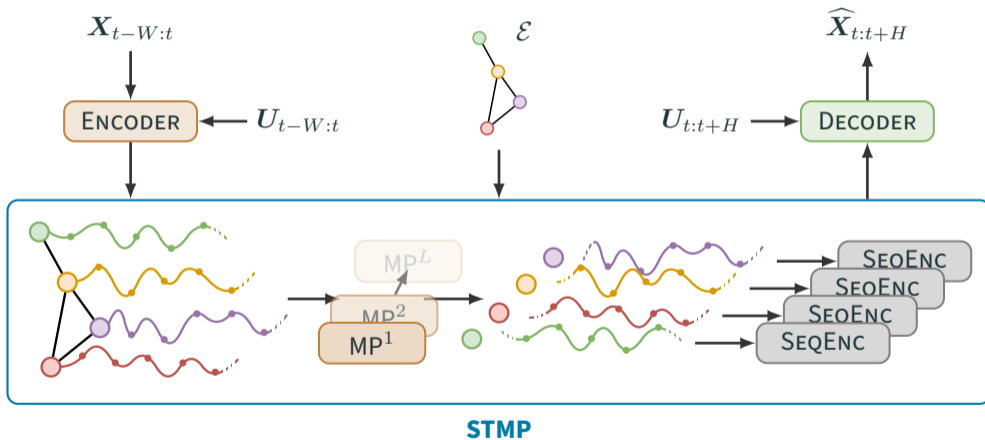
In STT approaches the two processing steps of TTS models are inverted:

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.

# Full STT model



# 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**.



Notebook

**Spatiotemporal Graph Neural Networks with tsl**

