

### Summer School 2025 CEA – EDF – INRIA

### From Graph Neural Networks to Dynamic Graphs: an introduction

Christophe Millet<sup>1,2</sup>

<sup>1</sup> ENS Paris-Saclay, 91190 Gif-sur-Yvette, France <sup>2</sup> CEA, DAM, DIF, 91297 Arpajon, France

Messier 74. Image released April 2, 2013. (Image credit: NASA/ESA/Hubble Heritage Team)



© MARK ANDERSON, WWW.ANDERTOONS.COM



"But as a proof of concept it's a total success."

### Motivation ••••• Graphs are everywhere

#### Graphs = vertices + edges

- Universal representation for structured data.
- Social networks (people, messages), biology (atoms/proteins, chemical bonds), sensor networks (stations, geodesic), recommendation (users/items, ratings), transportation (cities, roads/flights), etc.
- Graphs help us ask and sometimes answer fundamental questions about structure and interaction
  - What is the structure of a network ? Are there parts? (clustering)
  - Do the parts look the same? (similarity, isomorphism)
  - How can we model a set of graphs? (models)





[1] Kengkanna & Ohue, doi: 10.1038/s42004-024-0115-w, 2024.[2] Gaebler & Ceranna, doi:10.1007/s00024-020-02604-y, 2021.

### **Motivation** •••• When graphs are ignored

#### About the event

- Accidental explosion (Beirut, Aug. 4, 2020).
- 2.7 kt of Ammonium Nitrate (NH<sub>4</sub>NO<sub>3</sub>).

#### **Global to local graphs**

- Recorded by the IMS and regional networks (INSN, Israel National Seismic Network).
  - Infrasound stations:  $\sim 10^3$  km.
  - Seismic stations ( $\blacktriangle \land \land$ ): ~ 10<sup>2</sup> km.
- IS + seismic at  $10^2$  km (▲).
- Publicly available videos (~ 10<sup>2</sup> m).

### Magnitude/yield?

- Many magnitude estimates from empirical laws and various tech.  $\Rightarrow W \simeq 0.13 2$  kt TNT [\*].
- Solving the inverse wave propagation problem.



### **Motivation OOOO Does simulation help?**

### Yield estimates (Beirut, 2020)

 Using Green's functions of the wave equation (SEM3D) + signals.

 $\min_{W}\sum_{s=1}^{N} \|\mathbf{s}_{s}-\hat{\mathbf{s}}_{s}\|,$ 

- $\mathbf{s}_s = \mathbf{s}_0(W) *_t \mathbf{G}$ ;  $s_0$ : source model.
- $\hat{s}$ , s: recorded and simulated signals.
- Using a compressible flow solver and videos of shock dynamics.

260





### Motivation •••• GraphCast\* (2022)

#### An example of GNN

- 36.7 10<sup>6</sup> parameters.
- Multi-mesh derived from icosahedral meshes, from M<sup>0</sup> to M<sup>6</sup> (40962 nodes).
- Predicts ~ 100s of weather variables over 10 days at 0.25° resolution, in < 1 min.</li>
- Encoding/decoding with a single layer.
- **2** MPNN: L = 16 layers to propagate information from local to global scale.
  - \* arXiv 2212.12794



# **Preliminaries**Laplacians, GFT, convolution, spectral and spatial approaches

### 

#### Basic definitions

- Graph: G = (V, E), V: set of **nodes** (vertices);  $E \subset V \times V$ : set of **edges**.
- Edges can be directed or undirected.
- A graph is *connected* when, for any two nodes *i* and *j*, there exists a sequence of edges forming a path from *i* to *j*.
- Adjacency and degree matrices  $A, D \in \mathbb{R}^{|V| \times |V|}$

$$A_{ij} = \begin{cases} 1, \ (v_i, v_j) \in E\\ 0, \ (v_i, v_j) \notin E \end{cases} \text{ and } D_i = \sum_j A_{ij} \end{cases}$$

- For directed graphs, we have in-degree  $D_i^{\text{in}}$  and out-degree  $D_i^{\text{out}}$ .
- Learning tasks on graphs f: (G, X) → Y, X ∈ ℝ<sup>|V|×C</sup> (C: feature dim.).
   Node-level app. (graph clustering, classification of web pages).
   Edge-level app. (traffic prediction, recommendation systems).
   Graph-level app. (molecule classification, point cloud analysis).



- $\mathbf{A} = \begin{pmatrix} 0 & 1 & 0 \\ 1 & 0 & 1 \\ 0 & 1 & 0 \end{pmatrix}$
- $\mathbf{D} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 1 \end{pmatrix} \quad \mathbf{X} = \begin{pmatrix} 1 & 2 \\ 2 & 3 \\ 1 & 0 \end{pmatrix}$

### **Preliminaries** •••••••• **Combinatorial Laplacian**

#### Definition and key properties

Unnormalized Laplacian:

 $\mathbf{L} = \mathbf{D} - \mathbf{A},$ 

For 
$$\mathbf{X} \in \mathbb{R}^{|V|}$$
,  $(\mathbf{L}\mathbf{X})_i = \sum_j A_{ij}(x_i - x_j)$ .

Symmetric, positive semi-definite:

 $\mathbf{X}^{\mathrm{T}}\mathbf{L}\mathbf{X} = \frac{1}{2}\sum_{i}\sum_{j}A_{ij}(x_{i} - x_{j})^{2} \ge 0.$ 

 $\mathbf{X}^{\mathrm{T}}\mathbf{L}\mathbf{X} = 0$  (total variation = 0) characterizes signals with zero variation (see 2) across the graph, i.e.  $x_i = x_j$  for any (i, j).

Diagonalization: L = UΛU\*, with Λ = diag (λ<sub>1</sub>, ..., λ<sub>|V|</sub>), λ<sub>i</sub> ∈ ℝ.
 L1 = 0 ⇒ multiplicity m<sub>g</sub> of λ = 0 (dim(ker L)) equals the number of connected components (see 1).

An intuitive approach:  $\mathbf{U} = \text{nodes} \times \text{spectral modes}$ 

 $\mathbf{L} = \begin{pmatrix} \mathbf{Q} & \mathbf{0} \\ \mathbf{0} & \mathbf{Q} \end{pmatrix}; \ \mathbf{Q} = \begin{pmatrix} \mathbf{1} & -\mathbf{1} & \mathbf{0} \\ -\mathbf{1} & \mathbf{2} & -\mathbf{1} \\ \mathbf{0} & \mathbf{1} & \mathbf{1} \end{pmatrix}$  $\chi_{\mathbf{0}}(\lambda) = \lambda(\lambda - 1)(\lambda - 3)$  $\Rightarrow m_a(0) = \mathbf{2}$  $\mathbf{X}^{\mathrm{T}}\mathbf{L}\mathbf{X} = \frac{1}{2} \times 4 + \frac{1}{2} \times 16 = 10$ 

### **Preliminaries** 0000000 Other definitions

- Symmetric normalized Laplacian
  - L gives more weight to high-degree nodes (hubs) in the energy computation → topological bias.
  - Def.:  $L_{sym} = D^{-1/2}LD^{-1/2} = I D^{-1/2}AD^{-1/2}$ .
    - Counteract topol. bias by scaling **sender** and **receiver** contributions.
    - $\sigma(\mathbf{L}) \subset \left[0, \max_{i} D_{i}\right] \Rightarrow \sigma(\mathbf{L}_{sym}) \subseteq [0, 2], \text{ i.e. } \lambda_{max} \leq 2.$
  - Applications: spectral filtering, GCNs, Laplacian Encodings, ...

#### Random walk (rw) Laplacian

- Def.:  $L_{rw} = D^{-1}L = I D^{-1}A.$ 
  - Normalizes only the **sender**, consistent with a rw's view of the graph (each node distributes its information over its neighbors)  $\Rightarrow \sigma(\mathbf{L}_{rw}) \subseteq [0,2]$ .
  - Asym. operator that models **directional diffusion**, where information flows from  $v_i$  to its neighbors  $v_j$  with a **Markov transition matrix**  $\mathbf{P} = \mathbf{D}^{-1}\mathbf{A}$ .
- Applications: diffusion models, APPNP, Personalized PageRank, ...

 $\sigma = \{0, 1 \pm 1/\sqrt{2}\}$ 

 $P_{31} = 0$ 

 $P_{21} = 1$ 

 $\mathbf{L}_{\text{sym}} = \frac{1}{\sqrt{2}} \begin{pmatrix} \sqrt{2} & -1 & 0\\ -1 & \sqrt{2} & 0\\ 0 & 0 & 0 \end{pmatrix}$ 

 $\mathbf{L}_{\rm rw} = \frac{1}{2} \begin{pmatrix} 2 & -1 & -1 \\ -2 & 2 & 0 \\ 0 & 0 & 0 \end{pmatrix}$ 

**Fig.** Spectral and probabilistic interpretations of Laplacian for

a small directed graph. Rand.

walker located at  $v_i$  moves to

 $\mathbf{v}_i$  with  $P_{ii} = \Pr(i \rightarrow j)$ .

*P*<sub>12</sub>

### Preliminaries 0000000 A simple example (1/2)

- A (undirected) circular graph G = (V, E)
  - Each  $v_i$  is connected to  $v_{i+1}$  and  $v_{i-1}$  and

 $E = \{(i, i + 1 \mod |V|)\}.$ 

•  $\mathbf{D} = 2\mathbf{I}$  (each node is connected to 2 nodes),  $\mathbf{A} = \mathbf{J} + \mathbf{J}^{N-1}$ , where **J**: circular shift matrix of order N = |V|, i.e.

$$J_{ij} = \delta_{i,(j+1) \mod |V|}.$$

• For 
$$N = 4$$
:

$$\mathbf{J} = \begin{pmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 \end{pmatrix} \text{ and } \mathbf{J}^N = \mathbf{I}.$$

#### Use in energy-based models

- Circular graphs are used in energy-based models such as Convolutional or Temporal Restricted Boltzmann Machines.
- These architectures leverage circular connectivity to implement shared filters and translation invariance.

5 5 6 **Fig.** : Action of **A** on a circular graph. The operation aggregates forward and backward shifts of a signal, yielding a discrete circular convolution with kernel [1 0 1] and reflecting the graph's rotational invariance..





### Preliminaries 00000000 A simple example (2/2)

#### Spectral properties of the circular graph

Since eigenvalues of J are the *N*-th roots of unity  $(e^{2i\pi k/N})$ , the eigenpairs of L are  $(LU_k = \lambda_k U_k)$ :

$$A_k = 2[1 - \cos(2k\pi/N)] \in \mathbb{R}$$
  
 $(\mathbf{U}_k)_j = \frac{1}{\sqrt{N}} \exp\left(\frac{2k\pi i j}{N}\right).$ 

- Each  $\mathbf{U}_k \in \mathbb{C}^{|V|}$  (k = 0, ..., N 1) represents a harmonic mode over *G* and  $\lambda_k \in \mathbb{R}$  measures how much the mode oscillates.
  - All  $\lambda_k$  ( $k \neq 0, N/2$ ) appear in conjugate pairs, hence multiplicity 2 in the real spectrum (related to the graph's rotational symmetry).
- Signal reconstruction, for  $\mathbf{X} \in \mathbb{R}^{|V|}$ :

$$\mathbf{X} = \underbrace{\sum_{k=0}^{N-1} \hat{X}_k \mathbf{U}_k}_{\mathbf{U}\hat{\mathbf{X}}} \text{ with } \hat{X}_k = \mathbf{U}_k^* \mathbf{X},$$
  
Using  $\mathbf{U} = (\mathbf{U}_1 \cdots \mathbf{U}_{|V|})$  gives:  $\widehat{\mathbf{X}} = \mathbf{U}^* \mathbf{X}.$ 

 $\Rightarrow \widehat{\mathbf{X}}$  captures the signal's variation at different graph frequencies.



**Fig. :** Frequency modes for a circular graph with 8 nodes.

### **Preliminaries** 00000000 **GFT and convolution**

#### Definition of the Graph Fourier Transform (GFT)

• GFT allows to express  $\widehat{\mathbf{X}} \in \mathbb{R}^{|V|}$  in terms of modes on the graph:

GFT:  $\hat{\mathbf{X}} = \mathbf{U}^* \mathbf{X}$  and iGFT:  $\mathbf{X} = \mathbf{U} \hat{\mathbf{X}}$ .

- Total variation: or  $\mathbf{X}^{\mathrm{T}}\mathbf{L}\mathbf{X} = \sum_{k=0}^{N-1} \lambda_k |\hat{X}_k|^2$ .
- Parseval's identity:  $\mathbf{X}^{\mathrm{T}}\mathbf{X} = \widehat{\mathbf{X}}^*\widehat{\mathbf{X}} = \sum_{k=0}^{N-1} |\widehat{X}_k|^2$ .

### • Convolution on a graph *G*:

- No abelian group structure: unlike  $\mathbb{R}^n$  or  $\mathbb{Z}^n$  a general graph lacks translation invariance  $\Rightarrow$  no canonical convolution theorem.
- Convolution between  $\mathbf{f} \in \mathbb{R}^{|V|}$  and  $\mathbf{g} \in \mathbb{R}^{|V|}$

$$\mathbf{f} *_{\mathbf{G}} \mathbf{g} = \mathbf{U}(\underbrace{\widehat{\mathbf{f}} \odot \widehat{\mathbf{g}}}_{\text{diag}(\widehat{\mathbf{g}})\widehat{\mathbf{f}}}) \Rightarrow \mathbf{f} *_{\mathbf{G}} \mathbf{g} = \mathbf{U}\text{diag}(\widehat{\mathbf{g}})\mathbf{U}^*\mathbf{f}.$$

• **Spectral filtering**: we reinterpret  $\hat{\mathbf{g}}$  as the sampling of a function on the spectrum, i.e.  $\hat{g}_k = h(\lambda_k) \Rightarrow \mathbf{f} *_{\mathbf{G}} \mathbf{g} = \mathbf{U}h(\mathbf{\Lambda})\mathbf{U}^*\mathbf{f}$ .





### **Preliminaries** 00000000 Limitations: cospectrality

### **Spectral graph convolution**

Requires eigendecomposition of the graph Laplacian:

 $\mathbf{f} *_G \mathbf{g} = \mathbf{U}h(\mathbf{\Lambda})\mathbf{U}^*\mathbf{f}.$ 

- $\Rightarrow$  Impractical for large or changing graphs.
- High computational cost: requires full eigendecomposition of the Laplacian O(|V|<sup>3</sup>) → Lapack + machine-specific optimisations or Lanczos method for sparse networks.
- Graph-specific basis: eigenvectors U depend on graph topology ⇒ not transferable across graphs.
- Incompatible with dynamic graphs: any change in topology alters  $\sigma(\mathbf{A})$ ; same  $\sigma(\mathbf{A})$  does no imply same graph  $\Rightarrow$  cospectrality.



No locality: spectral filters are global by construction — they blend all node information through U.



**Fig.** : Fraction of **non-isomorphic tree** (connected and acyclic) pairs that share the same  $\sigma$ , as a function of graph size, for different matrix representations [1].

[1] Wilson & Zhu, Pattern Recognition, 2008.

### **Preliminaries** 0000000 From spectral to spatial theory

#### Polynomial approximation and locality

h can be approximated by a polynomial, leading to localized filters:

$$h(\lambda) \simeq \sum_{k=0}^{K} \theta_k \lambda^k \Rightarrow \mathbf{f} *_G \mathbf{g} = \sum_{k=0}^{K} \theta_k \mathbf{L}^k \mathbf{f},$$

where  $\mathbf{L}^k$  captures information from the *k*-hop neighborhood.





**Fig. :** Graph diffusion under powers of the Laplacian, starting from a Dirac signal at node 2

But this is still <u>fixed</u> (non-learnable structure, no adaptation to the task).

 Idea: replace global spectral filtering with local spatial message passing: information is propagated locally via learnable operations — a key idea behind Message Passing Neural Networks.

### **Graph Neural Networks** MPNN, GCN, over-smoothing/squashing, expressiveness, GAT, Graph Transformers



### 

#### GNN layer

• The basic idea is to update the state of each node through embedding in a  $C_{k+1}$ -dimensional space at each GNN layer:

$$\mathbb{R}^{|V| \times C_k}_{\mathbf{X}^{(k)}} \xrightarrow[layer k]{} \mathbb{R}^{|V| \times C_{k+1}}_{\mathbf{X}^{(k+1)}}$$

•  $\mathbf{x}_{j}^{(k)} \in \mathbb{R}^{C_{k}}$ : denotes the feature vector at node *j* (i.e., the *j*-th row of  $\mathbf{X}^{(k)}$ ).

■  $\mathbf{x}^{[i](k)} \in \mathbb{R}^{|V|}$ : denotes the *i*-th channel of the graph signal, i.e., the *i*-th column of  $\mathbf{X}^{(k)}$ .

$$\mathbf{X}^{[1]} = \begin{pmatrix} x_{11} & x_{12} \\ x_{21} & x_{22} \\ x_{31} & x_{32} \\ x_{41} & x_{42} \end{pmatrix} \mathbf{X}_{1}$$

#### Spatial vs spectral

- MPN (1): feature vector at each node is updated using messages from neighbors.
- GCN (2): each signal channel is transformed using graph convolution.





#### AGGregate-UPdate operation

The MP framework defines a GNN layer as:

$$\mathbf{h}_{j}^{(k)} = \mathbf{AGG}^{(k)}(\{\mathbf{x}_{i}^{(k)} | v_{i} \in \mathcal{N}(v_{j})\}),$$
$$\mathbf{x}_{j}^{(k+1)} = \mathbf{UP}^{(k)}\left(\mathbf{x}_{j}^{(k)}, \mathbf{h}_{j}^{(k)}\right).$$

- The AGGregate step can be either a fixed permutation-invariant operation (e.g., Σ<sub>j</sub>, avg, mean) or a learnable function.
- The UPdate step refines the aggregated message by combining it with the node's current state, allowing gating mechanisms (GRU) to control information flow.
- By adding self-loops in the graph:

$$\mathbf{x}_{j}^{(k+1)} = \sigma(\bigoplus_{i|v_{i} \in \mathcal{N}(v_{j})} f\left(\mathbf{x}_{i}^{(k)}, \mathbf{x}_{j}^{(k)}; \boldsymbol{\theta}^{(k)}\right))$$

- $\sigma$ : non-linear activation and  $\theta^{(k)}$  are trainable parameters.
- $\bigoplus$  is a symmetric operator:  $\sum$ , avg or max.



Google Research Blog. 2024

### **Graph neural networks** 00000000 Aggregate, Update and Pooling

- Choice of *f* 
  - *f* can be given by:

$$f(\mathbf{x}_i, \mathbf{x}_j; \mathbf{\theta}) = \mathbf{\theta} \mathbf{x}_i$$
  

$$f(\mathbf{x}_i, \mathbf{x}_j; \mathbf{\theta}) = \mathbf{\theta} (\mathbf{x}_i - \mathbf{x}_j)$$
 or MLP( $\mathbf{x}_i, \mathbf{x}_i - \mathbf{x}_j$ ).

- For  $\theta = I$  this resembles the action of the Laplacian:  $\mathbf{x}_{j}^{(k+1)} = -(\mathbf{L}\mathbf{X}^{(k)})_{j}$ .
- Aggregation over  $\kappa$ -hop neighbors introduces  $\kappa$  as a structural parameter. This controls the receptive field at each layer.
- Normalization with respect to  $\mathcal{N}(v_j)$  is often required to mitigate sensitivity to varying neighborhood sizes and to improve stability during training.
- **Graph pooling**: to obtain a graph-level representation  $\mathbf{y} \in \mathbb{R}^{C}$ :

$$\mathbf{y} = \bigoplus_{i \mid v_i \in V} \mathbf{x}_i,$$

⊕: denotes a permutation-invariant operation, used to compute a global graph representation for tasks like molecule classification or event detection.

• Other approaches: graph coarsening, learned pooling (DiffPool, SAGPool).

**Fig.** :  $\mathbf{x} \sim U([-1,1]), \bigoplus = \sum (+ \text{ norm.})$  3-hop neighborhood, with  $\mathbf{\theta} = \mathbf{I}$ . Neighborhood aggregation leads to **feature homogenization** at depth. Initial



June 19th, 2025 **19** 

### **Graph neural networks** 00000000 Graph Convolutional Networks (2)

#### Learning Convolutional Filters

Each channel CNN

$$\mathbf{x}^{[j](k+1)} = g(\sum_{i=1}^{C_k} \mathbf{x}^{[i](k)} *_G \mathbf{w}^{(k)}_{ij}),$$

where filters  $\mathbf{w}_{ij}^{(k)} \in \mathbb{R}^{|V| \times |V|}$  are approximated using Chebyshev polynomials, i.e.  $w(\lambda) = \sum_m \theta_m T_m(\lambda)$ 

$$\mathbf{x}^{[j](k+1)} = g(\sum_{i=1}^{C_k} \sum_{m=0}^{K} \theta_{ijm}^{(k)} \mathbf{T}_m(\mathbf{\tilde{L}}_{sym}) \mathbf{x}^{[i](k)}).$$

 Using T<sub>0</sub>(Ĩ<sub>Lsym</sub>) = I and T<sub>1</sub>(Ĩ<sub>Lsym</sub>) = Ĩ<sub>Lsym</sub> and keeping m = 0,1 gives the GCN approximation

$$\mathbf{x}^{[j](k+1)} = g(\sum_{i=1}^{C_k} \theta_{ij0}^{(k)} \underbrace{(\mathbf{I} + \mathbf{D}^{-\frac{1}{2}} \mathbf{A} \mathbf{D}^{-\frac{1}{2}})}_{\widetilde{\mathbf{A}}} \mathbf{x}^{[i](k)}).$$

In matrix form, the update rule  $h_{\theta^{(k)}}^{(k)} : \mathbb{R}^{|V| \times C_k} \to \mathbb{R}^{|V| \times C_{k+1}}$  is given by

$$\mathbf{X}^{(k+1)} = h_{\boldsymbol{\theta}^{(k)}}^{(k)} \big( \mathbf{X}^{(k)} \big) = g((\mathbf{I} + \mathbf{D}^{-\frac{1}{2}} \mathbf{A} \mathbf{D}^{-\frac{1}{2}}) \mathbf{X}^{(k)} \boldsymbol{\theta}^{(k)}).$$

 $\tilde{\mathbf{L}}_{sym} = 2\lambda_{max}^{-1}\mathbf{L}_{sym} - \mathbf{I} \Rightarrow \sigma(\tilde{\mathbf{L}}_{sym}) \subset [-1,1].$ Approx.  $\lambda_{max} = 2$  and  $\theta_{ij0}^{(k)} = -\theta_{ij1}^{(k)}.$ 



spectral locality  $\Rightarrow$  **spatial** locality

### **Graph neural networks**

00000000 MPNN and GCN drawbacks<sup>[1]</sup> Rusch, Bronstein, Mishra, arxiv abs/2303.10993, 2023. [2] Li, Han & Wu, AIAA Conf on AI, 2018.

[3] Xu et al., arxiv: abs/1806.03536, 2018.

4 layer(s)

**Fig.** : over-smoothing

 $X_{j}^{(0)} = \delta_{1j}, \mathbf{X}^{(k+1)} =$ 

 $\widetilde{\mathbf{A}}\mathbf{X}^{(k)}$ .

### • Over-smoothing [1,2]

• For a GCN:

$$\mathbf{X}^{(k+1)} = \sigma(\widetilde{\mathbf{A}}\mathbf{X}^{(k)}\mathbf{W}),$$

After *L* layers ( $\sigma = id$ ):

$$\mathbf{X}^{(K)} = \widetilde{\mathbf{A}}^K \mathbf{X}^{(0)} \prod_{k=0}^{K-1} \mathbf{W}^k.$$

 Repeated applications of the symmetrically normalized adjacency matrix A cause spectral filtering with decaying modes

 $\widetilde{\mathbf{A}}^{K}\mathbf{X}^{(0)} = \mathbf{U}\mathbf{\Lambda}^{K}\mathbf{U}^{*}\mathbf{X}^{(0)} \xrightarrow[K \to \infty]{} \mathbf{U}_{1}\mathbf{U}_{1}^{*}\mathbf{X}^{(0)}$ 

- Only the top eigenvector U<sub>1</sub> (often const.) survives (σ(Ã) ⊂ [−1,1]), leading to identical features for all nodes → over-smoothing.
- Influence distribution of the nodes [3]
  - Influence score  $I_u^{(k)}(v) = \text{Tr}(\partial_{\mathbf{x}_v^{(0)}} \mathbf{x}_u^{(k)}).$ 
    - For a standard MPNN or GCN,  $I_u^{(k)}(v) \sim \mathbb{E}(\mathbf{P}^k \mathbf{x}_u)$ .
    - This links MP to stochastic diffusion on the graph and explains the progressive loss of locality in deeper layers.

20 layer(s)

1 layer(s)

### **Graph neural networks** 00000000000 MPNN and GCN drawbacks

Over-squashing [1]

- Inability of GNNs to capture long-range dependencies due to compression of messages through paths γ of G.
- From update to paths ( $X \in \mathbb{R}^{|V| \times C}$ )

$$X_{j}^{(K)} = \sum_{j=1}^{|V|} \underbrace{\left(\widetilde{\mathbf{A}}_{j}^{K}\right)_{ij}}_{\sum_{\gamma \in \Gamma_{K}^{ij}} \prod_{(u,v) \in \gamma} \widetilde{\mathbf{A}}_{uv}} X_{j}^{(0)},$$

- $\Gamma_K^{ij}$ : set of all paths of length *K* from node *i* to node *j*.
- $\widetilde{\mathbf{A}}^{K}$ : aggregates the contributions of **all sequences of** *K* **hops** connecting *j* to *i*, each weighted by the product of edge weights along  $\gamma$ .
- Bottlenecks (Alon & Yahav, 2020)
  - When a large number of paths from distant nodes must pass through **a** limited set of intermediate nodes, these nodes are unable to carry sufficient information due to the fixed  $C \rightarrow$  over-squashing.



[1] Alon & Yahav, arxiv abs/2006.05205, 2020.



**Fig.** : A graph and multiple paths to connect  $v_i$  (i = 1, ..., |V|) and  $v_j$ .

### **Graph Neural Networks** 000000000 **GNNs expressiveness**

#### Permutation invariance

A model f(A, X) is permutation invariant if for all A ∈ ℝ<sup>|V|×|V|</sup>, all feature matrix X ∈ ℝ<sup>|V|×C</sup>, and all permutation matrix P ∈ ℝ<sup>|V|×|V|</sup>:
 f(PAP<sup>T</sup>, PX) = f(A, X)

Permutation equivariance:  $f(\mathbf{P}\mathbf{A}\mathbf{P}^{\mathrm{T}}, \mathbf{P}\mathbf{X}) = \mathbf{P}f(\mathbf{A}, \mathbf{X}).$ 

#### Graph isomorphism

• Two graphs  $G_1$  and  $G_2$  are isomorphic if  $\exists \mathbf{P} \in \mathbb{R}^{|V| \times |V|}$  s.t.  $\mathbf{P}\mathbf{A}_1\mathbf{P}^T = \mathbf{A}_2$ . Extension to tuples  $(G_1, \mathbf{X}_1)$  and  $(G_2, \mathbf{X}_2)$ :

 $\mathbf{P}\mathbf{A}_{1}\mathbf{P}^{\mathrm{T}} = \mathbf{A}_{2}$  and  $\mathbf{P}\mathbf{X}_{1} = \mathbf{X}_{2}$ .

- *f* permutation invariant and  $(G_1, \mathbf{X}_1) \sim (G_2, \mathbf{X}_2) \Rightarrow f(\mathbf{A}_1, \mathbf{X}_1) = f(\mathbf{A}_2, \mathbf{X}_2)$ .
- Graph isomorphism testing is NP-indeterminate problem [1].
- Heuristic: Weisfeler-Lehmann (1-WL) [2].
- Expressive power of many MPNN-based GNNs is <u>upper-bounded</u> by the 1-WL test — a heuristic for distinguishing non-isomorphic graphs.

[1] László Babai, arXiv:1512.03547, 2015.[2] Hamilton, Synthesis Lectures on AI & ML, 2022.



### **Graph Neural Networks** 00000000 **Graph Attention Networks**

#### Back to the aggregation step

The GAT framework [1] defines a GAT layer as:

$$\mathbf{x}_{j}^{(k+1)} = \sigma \left[ \sum_{i \mid v_{i} \in \mathcal{N}(v_{j})} \alpha_{ij} \mathbf{\Theta} \mathbf{x}_{i}^{(k)} \right],$$

or  $\mathbf{X}^{(k+1)} = \sigma(\mathbf{\alpha}^{\mathrm{T}} \mathbf{X}^{(k)} \mathbf{\theta}^{(k)\mathrm{T}})$ , where  $\alpha_{ij} = \operatorname{softmax}(e_{ij})$ 

$$e_{ij} = a(\mathbf{W}\mathbf{x}_i, \mathbf{W}\mathbf{x}_j)\mathbf{1}_{v_j \in \mathcal{N}(v_i)}.$$

- $a: \mathbb{R}^{C_l} \times \mathbb{R}^{C_l} \to \mathbb{R}$ : 1-layer MLP (+ LeakyReLU) or inner product/cos. similarity.
- Multi-head attention [2] introduces *M* coefficients  $\alpha_{ij}^m$  and  $\mathbf{x}_j^{(k+1)} = || \sigma(.)$ .
- Close to GCNs in form: the GAT layer can be seen as a GCN where  $\tilde{A}$  is replaced by a learnable attention matrix  $\alpha^{T}$ .

#### GAT and Beyond

- GAT layers remain permutation invariant (as GCNs).
- GATv2 [3] resolves expressiveness issues by allowing  $e_{ij} = \mathbf{a}^{\mathrm{T}} \sigma (\mathbf{W}[\mathbf{x}_i | | \mathbf{x}_j])$ .
- Limitations: like MPNNs, standard GATs are still limited by the 1-WL expressiveness [4].

[1] Velickovic et al., arxiv:1710.10903, 2017.
[2] Vaswani et al., NeurIPS, 2017.
[3] Brody et al., arxiv: 2105.14491, 2022.
[4] Xu et al., Arxiv:1810.00826, 2018

Fig. : attention-based message passing for  $v_i$ . Thicker arrows indicate higher importance.





### **Graph Neural Networks** 000000000 **Graph Transformers**

#### From local to global attention

■ Laplacian Positional Encoding [1,2],

for  $v_i$  we extract  $\mathbf{p}_j = (U_{i1} \cdots U_{ik}) \in \mathbb{R}^k$ 

- Serve as structural anchors for each node in the attention space, enabling the model to reason about graph topology.
- Allow the model to distinguish nodes with similar features but different topological roles, which is not possible using node features alone.
- Graph attention layers are defined by

$$\mathbf{h}_{j}^{(l+1)} = \left\| \begin{array}{c} {}_{k=1}^{K} \left[ \sum_{i \mid v_{i} \in \mathcal{N}(v_{j})} \alpha_{ij}^{lk} \mathbf{\theta}^{lk} \tilde{\mathbf{x}}^{(l)} \right], \\ \mathbf{x}_{j}^{(l+1)} = \operatorname{No} \circ f \circ \operatorname{No}(\mathbf{x}_{j}^{(l)} + \mathbf{h}_{j}^{(l)}), \end{array} \right.$$

- No: BatchNorm or LayerNorm; *f*: (nonlinear) MLP,  $\mathbf{x}_i^{(0)} = \mathbf{x}_i + \mathbf{W}\mathbf{p}_i$ ; attention coefficients  $\alpha_{ij}^{lk}$  are obtained from  $e_{ij}^{kl} = a(\mathbf{W}\mathbf{x}_i, \mathbf{W}\mathbf{x}_j)/\sqrt{C_l}$ .
- Global attention if  $\mathcal{N}(v_j) \to V$ .







## 

### **Graph Neural Networks**

### 

#### Add Features: enrich node/edge descriptors

- Add domain knowledge or structural information (degrees, distances, orbits, etc.).
- Random initial features to break symmetry [Sato et al., 2021].
- Subgraph-based features: counts, motifs, positional encodings  $\rightarrow$  **GSN**: Substruct. Net. [Bouritsas et al., 2022].
- Effective resistance & hitting times [Topping et al., 2022].
- Limitations: not always transferable; handcrafted choices may lack generality.

#### Modulate Message Passing: adapt interactions

- Use edge-dependent weights, attention, anisotropy  $\rightarrow$  GAT: Graph Attention Networks [Velickovic et al., 2018].
- Identity-aware aggregation [You et al., 2021].
- Directional aggregation from Laplacian eigenflows → DGN: Directional GNNs [Beaini et al., 2021].
- Can capture fine-grained node distinctions & directionality.
- **Tradeoff:** increased complexity, may amplify overfitting.

#### Modify the Graph: act on the computation structure

- Rewire the graph to enhance connectivity → DropEdge, digraph rewiring, adding virtual nodes.
- Use high-order structures  $\rightarrow$  k-WL GNNs [Morris et al., 2019], Ring-GNN [Chen et al., 2019].
- Subgraph GNN [Zhang et al., 2021].
- Encode hierarchy and locality → Nested GNNs, Hierarchical Pooling.
- **Drawback:** cost scales with higher-order terms.



## **3** Graph Neural Operators NO, GNO and Spatio-Spectral GNO



### **Graph Neural Operators**

\* Zongyi Li, Kovachki et al., 2021.

Input:  $[a(\mathbf{p}_i)]_i$ ,  $\mathbf{p} \in \mathbb{R}^3$ 

### Intuition of NO (t is ignored for simplicity)

- If G is the Green function of PDE  $\mathcal{L}(a)x = f$ , then x = G \* f.
- *G* is modelled as a kernel  $\kappa_{\theta}$ :  $G(\mathbf{p}, \mathbf{q}) \simeq \kappa_{\theta}(\mathbf{p}, \mathbf{q}, a(\mathbf{p}), a(\mathbf{q}))$  and for any  $\mathbf{x}: \mathbb{R}^d \to \mathbb{R}^m$ , a NO is defined by

$$\mathbf{K}_{\boldsymbol{\theta}}\mathbf{x}(\mathbf{p}) = \int \kappa_{\boldsymbol{\theta}}(\mathbf{p}, \mathbf{q}, a(\mathbf{p}), a(\mathbf{q})) \mathbf{x}(\mathbf{q}) d\mathbf{q}.$$

• FNO  $\kappa_{\theta}(\mathbf{p}, \mathbf{q}, c(\mathbf{p}), c(\mathbf{q})) = \kappa_{\theta}(\mathbf{p} - \mathbf{q}) \Rightarrow \mathbf{K}_{\theta}\mathbf{x}(\mathbf{p}) = \kappa_{\theta} * \mathbf{x}.$ 

#### Architecture of NOs

The mapping is learnt iteratively:

$$(a(\mathbf{p}), \mathbf{p}) \xrightarrow{P} \mathbf{X}^{(0)} \xrightarrow{F_1} \cdots \xrightarrow{F_L} \mathbf{X}^{(L)} \xrightarrow{Q} u(\mathbf{p}),$$
$$\mathbf{X}^{(l)} = \sigma_l \big( [\mathbf{W}^{(l)} + \mathbf{K}(a)] \mathbf{X}^{(l-1)} + \mathbf{b}^{(l)} \big).$$

- P: uplift layer, Q: projection layer.
- If K is a convolution kernel, the convolution theorem leads to:

$$\mathbf{K}(a)\mathbf{x} = \mathcal{F}^{-1}(\underbrace{\mathcal{F}(\mathbf{\kappa})}_{\mathbf{R}}\mathcal{F}(\mathbf{x}))$$

• The weights **R** are learnt inside each layer.



### **Graph Neural Operators •••** From FNOs to GNOs

#### Low-frequency truncation

- In practice R is truncated to the first m modes to obtain a low-rank operator acting on  $\widehat{\mathbf{X}}^{(\mathbf{k})}$ .
- Convolution in a FNO can be re-written as:

 $\mathbf{X}^{(k+1)} = \sigma(\underbrace{\mathcal{F}^{-1}}_{k} (\mathbf{R} \underbrace{\mathcal{F}(\mathbf{X}^{(k)})}_{k})))$ 

 $\mathbf{U}^{\mathrm{T}}\mathbf{X}^{(k)}$ Analogy GCN: U

•  $\mathbf{X}^{(k)} \in \mathbb{R}^{N \times C_k}$  multi-channel field,  $\mathbf{R} \in \mathbb{R}^{m \times m}$ .

#### Extension to graphs

The Fourier basis (FNO) can be replaced with the graph Laplacian eigenbasis U<sub>m</sub>:

 $\mathbf{X}^{(k+1)} = \mathbf{U}_m \mathbf{R} \times_1 \mathbf{U}_m^{\mathsf{T}} \mathbf{X}^{(k)}.$ 

- $\mathbf{U} \in \mathbb{R}^{|V| \times m}$ ,  $\mathbf{R} \in \mathbb{R}^{m \times C_k \times C_{k+1}}$  and  $\mathbf{X}^{(k)} \in \mathbb{R}^{|V| \times C_k}$ .
- $\times_1$  refers to a mode-1 tensor contraction.
- The MLP in a FNO can be replaced by a MPNN.

Fig. : 2D wave propagation in waveguides. **Input**: sound speed  $c: \mathbb{R}^2 \to \mathbb{R}$ , **output**: waveforms  $u: \mathbb{R}^2 \to \mathbb{R}$ ; 4 layers, 48 neurons/layer, m is varied.



29

### **Graph Neural Operators** ••• Spatio-Spectral GNO

- Global Spectral Graph Convolution
  - General form of graph convolution:

$$\mathbf{x}^{[j](k+1)} = \sigma(\sum_{i=1}^{C_k} \mathbf{x}^{[i](k)} *_G \mathbf{w}^{(k)}_{ij}),$$

- $*_G$ : Graph convolution, defined as filtering via  $\Lambda$ , i.e.  $\mathbf{x} *_G \mathbf{w} = \mathbf{U}g_{\theta}(\Lambda)\mathbf{U}^{\mathrm{T}}\mathbf{x}$ .
- $\mathbf{w}_{ij}^{(k)}$ : learnable spectral filters between channel *i* and *j*. For multi-valued  $\mathbf{X} \in \mathbb{R}^{|V| \times C_k}$ , we introduce  $\mathbf{R} \in \mathbb{R}^{|V| \times C_k \times C_{k+1}}$

$$\mathbf{X}^{(k+1)} = \sigma \left( \mathbf{U}_m \big( \mathbf{R} \times_1 (\mathbf{U}_m^{\mathrm{T}} \mathbf{X}^{(k)}) \big) \right).$$

- $U_m$ : first *m* eigenvectors of  $\Lambda$  computed using Locally Optimal Block Preconditioned Conjugate Gradient (LOBPCG) algorithm  $\rightarrow$  computational cost  $O(m\kappa|V|)$  if *G* is obtained using the  $\kappa$ -NN.
- Local Convolution with Spatial GNN
  - Aggregation is limited to the 1-hop direct neighbors  $\mathbf{x}_i^{(k)} \in \mathbb{R}^{C_k}$  $\mathbf{x}_i^{(k+1)} = \sum_{i \mid v_i \in \mathcal{N}(v_i)} \alpha_{ji} \mathbf{W} \mathbf{x}_i^{(k)}.$ 
    - $\alpha_{ij} = \sigma_2(\mathbf{W}_3 \sigma_1(\mathbf{W}_1[\mathbf{h}_j || \mathbf{h}_i || \mathbf{W}_2 w_{ij}]))$ : edge-weights are learned within a gating mechanism, with  $w_{ij} = ||\mathbf{p}_i \mathbf{p}_j||$ , where  $\mathbf{p}_i$  is the positions of node *i* and  $\mathbf{h}_i$  is its Lipschitz embedding\*.

\* Sarkar & Chakraborty, 2409.00604, 2024.



**Fig.** : Neural architecture of a block. It exploits spectral and spatial GNN to formulate the kernel integration operator [1].

*f* concatenates the outputs of Spectral and Spatial GNNs.

### **Dynamic GNN for seismology** Architectures, performance, interpretability, over-smoothing, <u>practice</u>

### **Dynamic GNN for Seismology**

### 

### Goal of the approach:

Learn a parametric model

• Input:  $|V| \times d$  signals recorded by |V| stations and sampled using *P* time steps. *d*: NS, EW, Up-Down components.

 $\mathbf{S} \mapsto \mathbf{y}$ 

 $f_{\mathbf{A}} : \overset{\bullet}{\mathbb{S}} \to \mathbb{R}^4$ 

- Output: characteristics of events.
- Loss function:

$$\mathcal{L} = \frac{1}{4n} \sum_{k=1}^{n} \sum_{j=1}^{4} |\hat{y}_{kj} - y_{kj}|.$$

•  $y_{kj}$ : latitude, longitude, depth, or magnitude of event k.

### **Graph-related issues**

- The graph is not specified and A is *a priori* not known.

Fig. Southern California Seismic Network



### 

**1** Convolutional Encoding  $h_{\mathbf{q}^{(I)}}^{(I)} \colon \mathbb{R}^{3C_0} \mapsto \mathbb{R}^{d^{(I)}}$ 

Extraction of temporal features (CNN-based encoder):

$$x_{s}^{[i](l)}[n] = \operatorname{ReLu}(b_{i}^{(l)} + \sum_{j=1}^{d_{l-1}} \underbrace{\mathbf{W}_{ji}^{(l)} * x_{s}^{[j](l-1)}[n]}_{\sum_{k=0}^{4} w_{jik}^{(l)} x_{s}^{[j](l-1)}[n+k]}).$$

- $x_s^{[i](l)}[n]$ : value of the *i*-th feature channel at time index *n* for station *s* at layer *l*.
- Valid output indices:  $n < m_{l-1} 4 \Rightarrow$  temporal dimension shrinks accross layers.
- 2 Spatial information integration  $h_{\mathbf{\theta}^{(\mathrm{II})}}^{(\mathrm{II})} : \mathbb{R}^{(d^{(\mathrm{I})}+2) \times |V|} \to \mathbb{R}^{d^{(\mathrm{II})} \times |V|}$ 3 Aggregation Pool:  $\mathbb{R}^{d^{(\mathrm{II})} \times |V|} \to \mathbb{R}^{d^{(\mathrm{II})}}$
- **4** Prediction  $h_{\boldsymbol{\theta}^{(\mathrm{III})}}^{(\mathrm{III})} : \mathbb{R}^{d^{(\mathrm{II})}} \to \mathbb{R}^4$
- Extraction of graph-level information:

$$\mathbf{y} = \sigma_2^{(\mathrm{III})} (\mathbf{W}_2^{(\mathrm{III})} \sigma_1^{(\mathrm{III})} \left( \mathbf{W}_1^{(\mathrm{III})} \mathbf{z} + \mathbf{b}_1^{(\mathrm{III})} \right) + \mathbf{b}_2^{(\mathrm{III})})$$

■  $\mathbf{z} \in \mathbb{R}^{d^{(II)}}$  is derived by aggregating the features  $\mathbf{X} \in \mathbb{R}^{d^{(II) \times |V|}}$  using Pool( $h_{\theta^{(II)}}^{(II)}(.)$ ).



 $\mathbf{s}_{s}^{\left[ \cdot \right]} \in \mathbb{R}^{d \times P}$ 



### **Dynamic GNN for Seismology** 00000000000 Spatio-Temporal Graph Local Convolution

### ST-GC Block (based on MP $\rightarrow$ /ocal)

Geographic distance

$$\mathbf{x}_{s,\text{geo}}^{(l+1)} = \max_{j \mid v_j \in \mathcal{N}_{\text{geo}}(v_s)} f_{\text{geo}}^{(l)}(\mathbf{u}_s^{(,l)}, \mathbf{u}_j^{(l)} - \mathbf{u}_s^{(l)}),$$

- $\mathbf{u}_s = \mathbf{x}_s || \mathbf{p}_s \ (\mathbf{p}_s \in \mathbb{R}^2)$  allow the model to distinguish nodes with similar features but different topological roles.
- Asymmetric update  $\mathbb{R}^{d_{\text{geo}}} \to \mathbb{R}^{d^{(1)}}$ ;  $d^{(1)}$ : dimension of the encoded features; input dim.:  $2(d^{(1)} + 2)$ .
- $\mathcal{N}_{geo}^{(l)}$ :  $\kappa$ -NN of based on geographic positions  $\|\mathbf{p}_i \mathbf{p}_j\|_2$ .
- Feature similarity

$$\mathbf{x}_{s,\text{dist}}^{(l+1)} = \max_{j \mid v_j \in \mathcal{N}_{\text{sig}}(v_s)} f_{\text{sig}}^{(l)}(\mathbf{x}_s^{(,l)}, \mathbf{x}_j^{(l)} - \mathbf{x}_s^{(l)}),$$

- Asymmetric update  $\mathbb{R}^{d_{sig}} \rightarrow \mathbb{R}^{d^{(1)}}$ ; input dim.:  $2d^{(1)}$ .
- $\mathcal{N}_{sig}^{(l)}$ :  $\kappa$ -NN with  $\|\mathbf{x}_i \mathbf{x}_j\|^2$ : combine  $\mathbf{x}_i$ 's from distant stations.
- Multi-scale embedding  $f: \mathbb{R}^{(L+1)d^{(1)} \times |V|} \to \mathbb{R}^{d^{(1)} \times |V|}$ 
  - $\mathbf{X}^{(0)} \mid\mid \mathbf{X}^{(1)} \mid\mid \mathbf{X}^{(2)} \cdots \mid\mid \mathbf{X}^{(L)}$  where  $\mathbf{X}^{(l)} = \mathbf{X}^{(l)}_{\text{geo}} \oplus \mathbf{X}^{(l)}_{\text{dist}}$





**Fig.** : ST-GC Block and multi-scale embedding.

## **Dynamic GNN for Seismology**

#### **Convolutional encoding (shared)**

- 12 convolutional layers (kernel size = 5).
- Max pooling applied every 3 layers (kernel size = 4).

### **Spatial information integration**

- Edgeless-GNN and Spectral-GNN
  - 2-layer MLP followed by max-pooling.
  - Hidden layer size:  $2^7 = 128 \Rightarrow d^{(II)} = 128$ .
  - For (local) GCN,  $\kappa = 5$  spatial neigbors per node.

#### Dynamic-GNN

- ST-GC: 2-layer MLPs (one MLP per graph: geometric & feature-based). ST-GC block is repeated L = 4 times.
- Hidden layer size: 32 per branch.
- $\kappa = 5, L = 4 \text{ layers} \Rightarrow d^{(\text{III})} = (L+1)d^{(\text{I})} = 5 \times 2^6 = 320.$

### **Prediction**

2-layer MLP with a hidden size of 128.



#### 0000000000000 Datasets



- 20 yrs (01/01/2000 06/30/2019).
- 3493 events with  $M_l > 2$ .
- **72** stations, **3** components.

SCSN<sup>2</sup> (USA)





#### Variance decomposition

$$\sigma^{2} = \frac{\frac{1}{KR} \sum_{k=1}^{K} \sum_{r=1}^{R} (\bar{E}_{k} - \bar{E}_{k})^{2}}{\sigma_{\text{opt}}^{2}} + \frac{\frac{1}{K} \sum_{k=1}^{K} (\bar{E}_{k} - \bar{E})^{2}}{\sigma_{\text{dat}}^{2}}.$$

- $\sigma_{\rm opt}^2$ : seed-level variability.
- $\sigma_{dat}^2$ : variability induced by differences between test folds.

### Main results Edgeless vs Dynamic

- Observations
  - $\sigma^2$  is primarly driven by  $\sigma_{dat}^2$  as  $M_w \nearrow$ .
  - $\sigma_{dat}^2$  > with Dynamic-GNN.

### Explanations

- Architecture: Dynamic-GNN leverages adaptive topology.
- **Distribution**: More data reduces average RMSE, but imbalanced  $M_w$  distribution  $\Rightarrow \sigma^2 \nearrow$  for rare (large) events.
- **Optimization**: seed-dependent variability accounts for 1/3 of  $\sigma^2$ ; mainly affects low- $M_w$  events with weaker feature signals and higher SNR.



0.02

0.01

0.00 +<del>-</del> 2.5

**Dynamic-GNN** 

3.0

3.5

### **Applications for seismology** 0000000000 **Robustness to input signals**

#### **Three strategies**

 $\mathbf{X}_{S0} \rightarrow \mathbf{X}_{S1} \rightarrow \mathbf{X}_{S2}$ 

- S0 (standard): reference configuration
  - Broadband filtering (1–8 Hz): removes LFs below the corner frequency (~ 0.2 Hz), which are critical for seismic moment estimation.
  - Aligned time windows, based on origin time  $t_0$  of each event.
- S1 (narrowband): filtered waveforms
  - Apply a bandpass filter (1–1.5 Hz).
  - Removes much of the spectral content related to rupture complexity (e.g., HF scattering outside the band).
- S2 (random sliding): **misaligned** windows
  - Apply a random shift  $(10^2 U([-1,1]))$  to the start time of each window.
  - Add station-specific time shifts compare to published works\*.
  - Breaks alignment to theoretical arrival times, simulating poor or missing phase picks (e.g., early warning, low-SNR environments).





**00000000000000 Gains depend on graph and data quality** 





### **Trends**

- Why doesn't Dynamic-GNN outperform Edgeless-GNN on Resif<sub>2.5</sub> and S2? (1)
- S2 do not impact localization performance (2)
  - A proxy scale can be defined as:

 $\lambda = \frac{1}{N(\kappa-1)} \sum_{i=1}^{N} \sum_{j=1}^{\kappa-1} d_{\text{geo}}(v_i, n_j(v_i)).$ 

 $d_{\text{geo}}$ : geodesic distance to the neighbors  $n_j(v_i)$ 

- $\lambda \simeq 104$  km for Resif<sub>2.5</sub> and  $\lambda \simeq 39$  km for SCSN.
- GNN are able to resolve **sub-\lambda spatial** differences.

- The retained events (*M<sub>w</sub>* > 2.5) are less sensitive to SNR effects, making GNN models equally robust.
- The Resif<sub>2.5</sub> network is sparse and uneven, reducing the ability of dynamic graphs to extract meaningful spatial patterns.
- The limited number of training events (≈1000) may not support the higher model complexity of Dynamic-GNN.

### **Sensitivity indices**

• Contribution of node *s* to the final prediction:

$$\left[S_{s}^{(L)}\right]^{2} = \left\|\nabla_{\mathbf{x}_{s}^{(l)}}\epsilon\right\|_{2}^{2}.$$

Using  $\mathbf{z} = \max_{s} \mathbf{x}_{s}$  leads to:

$$\left[S_{s}^{(L)}\right]^{2} = \sum_{j=1}^{d_{z}} \mathbb{1}_{s=s_{j}^{(L)}} \left(\partial_{z_{j}} \epsilon\right)^{2}.$$

 $S_s^{(L)}$  acts as a station-wise attribution **score**, highlighting which nodes contributed most to the prediction error  $\epsilon$ .

### **Layer-Dependent Spatial Attention**

- Deeper GNN layers focus on fewer, highsensitivity stations  $V^{(l)} \subset V$ :  $|V^{(l)}|/|V| \rightarrow h.$
- Stations cluster near the epicenter, resembling a learned spatial attention mechanism.



### **Oversmoothing**

- Deeper GNNs tend to oversmooth, i.e. make node embeddings indistinguishable across G<sub>geo</sub> and G<sub>sig</sub>.
- For GCNs with fixed topology and linear propagation, *∧* L causes the node embeddings to converge towards a low-dim. invariant subspace: d<sub>M</sub>(X<sup>(l)</sup>) ≤ (sλ)<sup>l</sup>d<sub>M</sub>(X<sup>(0)</sup>) [1].

### Dirichlet energy $Tr(\mathbf{X}^T \mathbf{L}_{sym} \mathbf{X})$ [2]

Measures local variation of node embeddings (low *E* = smoothness). For directed graphs:

$$\mathcal{E}(\mathbf{X}^{(l)}) = \sum_{i} \sum_{j \mid v_j \in \mathcal{N}(v_i)} \left\| \mathbf{e}_{ij}^{(l)} \right\|_2^2$$

- $\mathbf{e}_{ij}^{(l)} = \mathbf{x}_i^{(l)} (d_i^{\text{in}})^{-1/2} \mathbf{x}_j^{(l)} (d_j^{\text{out}})^{-1/2}$  with the dimensions defined as  $d_i^{\text{in}} = \kappa$  and  $d_j^{\text{out}} = \#\{v_i | v_j \in \mathcal{N}(v_i)\}$ .
- This double normalization ensures  $\mathcal{E}$  remains well-defined when  $\mathcal{N}(v_i)$ 's are defined via attention mechanisms.

[1] Oono & Suzuki, arxiv 1905.10947, 2019 [2] Rusch, Bronstein & Mishra, arxiv 2303.10993, 2023



**Fig.** Averaged Dirichlet energies over all the events and  $\pm \sigma$ . Exponential decrease is satisfied which confirms oversmoothing. Resif<sub>2.5</sub>.



### Concat in **2** mitigates oversmoothing

- $\kappa = |\mathcal{N}(v_i)|$  defines the local receptive field (= 5)
  - Stacking multiple layers means that the local neighborhoods are recursively propagated across layers.
  - $\kappa$  and *L* should be <u>co-tuned</u> to avoid oversmoothing.
- || provides multi-scale agreggation (③) that maps  $\widetilde{\mathbf{X}}^{(L)} \in \mathbb{R}^{Ld}$  into a lower dimensional space

$$\underbrace{\mathbf{X}^{(0)} || \cdots || \mathbf{X}^{(L)}}_{\widetilde{\mathbf{X}}^{(L)}} \mapsto z = \operatorname{Pool}(\widetilde{\mathbf{X}}^{(L)}) \in \mathbb{R}^{d}.$$

- $X^{(0)} = X \parallel P$  (X: output of 1) inspired by Positional Encodings.
- I Plays a similar role to multi-head attention in Graph Transformers: aggregating multi-scale relational patterns into a unified representation
- || can be applied to large classes of GNN.
  - For Edgeless\*-GNN and "poor" datasets (**Resif**<sub>2.5</sub>)  $\rightarrow$  performance.
  - Preserve express. for GCNs (Chen *et al.*, arxiv 2007.02133, 2020).

**Fig.** Prediction error with *Dynamic*-GNN.  $\overline{\text{MAE}} \pm \sigma$  (solid lines);  $\overline{\text{RMSE}}$  (dashed line) **Resif**<sub>2.5</sub>.

		/
Туре	$\overline{E}\pm\sigma_{ m total}$	- //
RMSE	$\textbf{0.175} \pm \textbf{0.019}$	
MAE	$\textbf{0.139} \pm \textbf{0.018}$	Ref*
RMSE	$0.132\pm0.015$	
MAE	$0.101\pm0.013$	
RMSE	$\textbf{0.129} \pm \textbf{0.013}$	
MAE	$\textbf{0.098} \pm \textbf{0.012}$	
	€ <sub>Mw</sub> 0.1	without
	•	2 6 12 20 # layers ( <i>L</i> )

#### **Dynamic Graphs: Physics-aware and interpretable**

- The Dynamic-GNN learns adaptive station-to-station connectivity, reflecting underlying physical processes like wavefront propagation.
- Sensitivity scores reveal how information is routed, resembling a spatial attention mechanism focusing on informative stations near the epicenter.

#### **Concatenation mitigates oversmoothing**

 Aggregating intermediate outputs provides multi-scale information and preserves node-wise diversity (conceptually aligned with multi-head attention in Transformers, where different attention heads specialize in different relational patterns).

#### **Robustness in low-data settings**

- The Dynamic-GNN adapts its graph to the signal structure, improving generalization on small or imbalanced datasets (Resif<sub>2.5</sub> vs SCSN).
- Fixed-topology models struggle in such settings, as they cannot exploit dataspecific spatial correlations.



## **Open questions**

#### **Over-SMoothing (OSM)**

#### OSM is not always harmful

- OSM is not systematically detrimental to performance it can help classification if it aligns with the task objective (Fig.).
- Keriven et al. (2022) suggest: some OSM is desired, especially in homophilic graphs.

#### OSM happens non-uniformly

- Smoothing can occur faster in some embedding subspaces than others. If labels are aligned with these slow-smoothing subspaces → performance improves.
- This leads to the idea of task-oriented smoothing.

#### OSM analysis in real-world GNNs

- Real GNNs may not suffer as much from OSM, thanks to residual connections, gated or relational variants, etc.
- There is a need to extend OSM analysis to real-world architectures and not idealized deep GCNs.



**Fig.** Evolution of node embeddings from a GCN under increasing message-passing depth, showing class separation collapse as smoothing increases.

## **Open questions**

#### **Over-SQuashing (OSQ)**

### Compression vs. Bottlenecks

#### **Fig.** Rewiring the graph improves connectivity by reducing bottlenecks. The transformation $\mathcal{R}(G)$ increases $\lambda_2$ , lowers resistance $R_{tot}$ , and facilitates long-range message passing.

 $\lambda_2 = 0.05$  and  $R_{tot} = 819$ 





- OSQ arises from exponential compression of messages from many paths through a limited number of edges (bottlenecks).
- Classical metrics  $(\lambda_2, R_{tot})$  often used, but Structural metrics may fail to capture task-relevant information flow, especially when node labels or signals are misaligned with structural bottlenecks.

#### Bottleneck metrics vs. task sensitivity

- Metrics like curvature, resistance distance identify structural bottlenecks.
- But they might not align with what the network is trying to learn → There's a **gap between structural measures and task relevance**.

### Spectral perspective

- Spectral decomposition offers a way to measure how label information is distributed across frequency modes:
  - Homophilic labels: aligned with **low-frequency** modes.
  - Heterophilic labels: encoded in high-frequency components.





# Any questions now or later?

Christophe.millet@cea.fr



Abboud, R., Dimitrov, R., & Ceylan, I. I. (2022, December). Shortest path networks for graph property prediction.
 In Learning on Graphs Conference (pp. 5-1). PMLR.

•Abu-EI-Haija, Sami, et al. "Mixhop: Higher-order graph convolutional architectures via sparsified neighborhood mixing." international conference on machine learning. PMLR, 2019.

•Alev, V. L., Anari, N., Lau, L. C., & Gharan, S. O. (2018). Graph clustering using effective resistance http://drops.dagstuhl.de/opus/volltexte/2018/8369.

Alon, U., & Yahav, E. (2020). On the bottleneck of graph neural networks and its practical implications. In ICLR 2021.

 Arnaiz-Rodríguez, A., Begga, A., Escolano, F., & Oliver, N. (2022). Diffwire: Inductive graph rewiring via the Lovász bound. In the First Learning on graphs (LoG) Conference 2022.

•Arnaiz-Rodriguez, A., Curto, G., & Oliver, N. (2024). Structural Group Unfairness: Measurement and Mitigation by means of the Effective Resistance. In TrustLOG Workshop at WWW 2024.

Banerjee, P. K., Karhadkar, K., Wang, Y. G., Alon, U., & Montúfar, G. (2022, September). Oversquashing in gnns through the lens of information contraction and graph expansion. In 2022 58th Annual Allerton Conference on Communication, Control, and Computing (Allerton) (pp. 1-8). IEEE.

•Chamberlain, B., Rowbottom, J., Eynard, D., Di Giovanni, F., Dong, X., & Bronstein, M. (2021). Beltrami flow and neural diffusion on graphs. Advances in Neural Information Processing Systems, 34, 1594-1609.

•Bi, W., Du, L., Fu, Q., Wang, Y., Han, S., & Zhang, D. (2022). Make heterophily graphs better fit gnn: A graph rewiring approach. arXiv preprint arXiv:2209.08264.

•Black, M., Wan, Z., Nayyeri, A., & Wang, Y. (2023, July). Understanding oversquashing in gnns through the lens of effective resistance. In International Conference on Machine Learning (pp. 2528-2547). PMLR.

- •Bodnar, C., Di Giovanni, F., Chamberlain, B., Lio, P., & Bronstein, M. (2022). Neural sheaf diffusion: A topological perspective on heterophily and oversmoothing in gnns. Advances in Neural Information Processing Systems, 35, 18527-18541.
- •Brüel-Gabrielsson, R., Yurochkin, M., & Solomon, J. (2022). Rewiring with positional encodings for graph neural networks. TMLR 2023
- Buchnik, E., & Cohen, E. (2018, June). Bootstrapped graph diffusions: Exposing the power of nonlinearity. In Abstracts of the 2018 ACM International Conference on Measurement and Modeling of Computer Systems (pp. 8-10).
- •Cai, C., & Wang, Y. (2020). A note on over-smoothing for graph neural networks. In ICML 2020.
- •Cai, C., Hy, T. S., Yu, R., & Wang, Y. (2023, July). On the connection between mpnn and graph transformer. In International Conference on Machine Learning (pp. 3408-3430). PMLR.
- •Chamberlain, B., Rowbottom, J., Gorinova, M. I., Bronstein, M., Webb, S., & Rossi, E. (2021, July). Grand: Graph neural diffusion. In International conference on machine learning (pp. 1407-1418). PMLR.
- •Chamberlain, B., Rowbottom, J., Eynard, D., Di Giovanni, F., Dong, X., & Bronstein, M. (2021b). Beltrami flow and neural diffusion on graphs. Advances in Neural Information Processing Systems, 34, 1594-1609.
- •Chen, M., Wei, Z., Huang, Z., Ding, B., & Li, Y. (2020, November). Simple and deep graph convolutional networks. In International conference on machine learning (pp. 1725-1735). PMLR.
- •Chen, T., Zhou, K., Duan, K., Zheng, W., Wang, P., Hu, X., & Wang, Z. (2022). Bag of tricks for training deeper graph neural networks: A comprehensive benchmark study. IEEE TPAMI.
- Chung, F. R. (1997). Spectral graph theory (Vol. 92). American Mathematical Soc..
- •Devriendt, K., & Lambiotte, R. (2022). Discrete curvature on graphs from the effective resistance. Journal of Physics: Complexity, 3(2), 025008.

- Di Giovanni, F., Giusti, L., Barbero, F., Luise, G., Lio, P., & Bronstein, M. M. (2023, July). On over-squashing in message passing neural networks: The impact of width, depth, and topology. ICML
- Di Giovanni, F., Rowbottom, J., Chamberlain, B. P., Markovich, T., & Bronstein, M. M. (2023). Understanding convolution on graphs via energies. In TLMR.
- Di Giovanni, F., Rusch, T. K., Bronstein, M. M., Deac, A., Lackenby, M., Mishra, S., & Veličković, P. (2024). How does over-squashing affect the power of GNNs?. TMLR.
- Dwivedi, V. P., Rampášek, L., Galkin, M., Parviz, A., Wolf, G., Luu, A. T., & Beaini, D. (2022). Long range graph benchmark. Advances in Neural Information Processing Systems, 35, 22326-22340.
- •Eliasof, Moshe, Eldad Haber, and Eran Treister. "Pde-gcn: Novel architectures for graph neural networks motivated by partial differential equations." Advances in neural information processing systems 34 (2021): 3836-3849.
- •Errica, F., Christiansen, H., Zaverkin, V., Maruyama, T., Niepert, M., & Alesiani, F. (2024). Adaptive Message Passing: A General Framework to Mitigate Oversmoothing, Oversquashing, and Underreaching. In JMLR 2024
- •Fesser, L., & Weber, M. (2024, April). Mitigating over-smoothing and over-squashing using augmentations of Forman-Ricci curvature. In Learning on Graphs Conference (pp. 19-1). PMLR.
- •Gasteiger, J., Weißenberger, S., & Günnemann, S. (2019). Diffusion improves graph learning. Advances in neural information processing systems, 32.
- •Giraldo, J. H., Skianis, K., Bouwmans, T., & Malliaros, F. D. (2023, October). On the trade-off between oversmoothing and over-squashing in deep graph neural networks. In ICKM.
- •Gutteridge, B., Dong, X., Bronstein, M. M., & Di Giovanni, F. (2023, July). Drew: Dynamically rewired message passing with delay. In International Conference on Machine Learning (pp. 12252-12267). PMLR.

•Hamilton, W. L. (2020). Graph representation learning. Morgan & Claypool Publishers.

•Hasanzadeh, A., Hajiramezanali, E., Boluki, S., Zhou, M., Duffield, N., Narayanan, K., & Qian, X. (2020, November). Bayesian graph neural networks with adaptive connection sampling. In ICML 2022.

•Huang, K., Wang, Y. G., & Li, M. (2024). How Universal Polynomial Bases Enhance Spectral Graph Neural Networks: Heterophily, Over-smoothing, and Over-squashing. arXiv preprint arXiv:2405.12474.

-Jamadandi, Adarsh, Celia Rubio-Madrigal, and Rebekka Burkholz. "Spectral Graph Pruning Against Over-Squashing and Over-Smoothing." arXiv preprint arXiv:2404.04612 (2024).

 Karhadkar, K., Banerjee, P. K., & Montúfar, G. (2022). FoSR: First-order spectral rewiring for addressing oversquashing in GNNs. In ICLR 2023

•Jiang, W., Liu, H., & Xiong, H. (2023). Survey on Trustworthy Graph Neural Networks: From A Causal Perspective. arXiv preprint arXiv:2312.12477.

•Keriven, N. Not too little, not too much: a theoretical analysis of graph (over) smoothing. NeurIPS 2022.

•Kondor, R. I., & Lafferty, J. (2002, July). Diffusion kernels on graphs and other discrete structures. In Proceedings of the 19th international conference on machine learning (Vol. 2002, pp. 315-322).

•Li, Q., Han, Z., & Wu, X. M. (2018, April). Deeper insights into graph convolutional networks for semi-supervised learning. In Proceedings of the AAAI conference on artificial intelligence (Vol. 32, No. 1).

•Li, G., Muller, M., Thabet, A., & Ghanem, B. (2019). Deepgcns: Can gcns go as deep as cnns?. In Proceedings of the IEEE/CVF international conference on computer vision (pp. 9267-9276).

•Liu, M., Gao, H., & Ji, S. (2020, August). Towards deeper graph neural networks. In Proceedings of the 26th ACM SIGKDD international conference on knowledge discovery & data mining (pp. 338-348).

Liu, Y., Zhou, C., Pan, S., Wu, J., Li, Z., Chen, H., & Zhang, P. (2023, April). Curvdrop: A ricci curvature based approach to prevent graph neural networks from over-smoothing and over-squashing. In WWW.

•Liu, Y., Zheng, Y., Zhang, D., Lee, V. C., & Pan, S. (2023b, June). Beyond smoothing: Unsupervised graph representation learning with edge heterophily discriminating. AAAI.

- Lim, D., et al. "New benchmarks for learning on non-homophilous graphs". In WWW Workshop on GLB, 2021.
- Lovász, L. (1993). Random walks on graphs. Combinatorics, Paul erdos is eighty, 2(1-46), 4.
- Luan, S., Hua, C., Lu, Q., Zhu, J., Zhao, M., Zhang, S., ... & Precup, D. (2022). Revisiting heterophily for graph neural networks. Advances in neural information processing systems, 35, 1362-1375.
- Luan, S. et al (2024). The Heterophilic Graph Learning Handbook: Benchmarks, Models, Theoretical Analysis, Applications and Challenges. Arxiv 2407.09618
- •Ma, Y., Liu, X., Shah, N., & Tang, J. (2022). Is homophily a necessity for graph neural networks?. ICLR 2022.
- Maskey, S., Paolino, R., Bacho, A., & Kutyniok, G. (2024). A fractional graph laplacian approach to oversmoothing.
   Advances in Neural Information Processing Systems, 36.
- •Newman, M. "Assortative mixing in networks". Phys. Rev. Lett., 89, 2002.
- •Nguyen, K., Hieu, N. M., Nguyen, V. D., Ho, N., Osher, S., & Nguyen, T. M. (2023, July). Revisiting over-smoothing and over-squashing using ollivier-ricci curvature. In ICML 2023.
- Oono, K., & Suzuki, T. (2019). Graph neural networks exponentially lose expressive power for node classification. In ICLR 2020.
- Pei, H. et al. "Geom-GCN: Geometric GCNs". In ICLR, 2019.

cea

Pham, T., Tran, T., Dam, H., & Venkatesh, S. (2017). Graph classification via deep learning with virtual nodes. arXiv preprint arXiv:1708.04357.

•Qiu, H., & Hancock, E. R. (2006). Graph embedding using commute time. In Structural, Syntactic, and Statistical Pattern Recognition: Joint IAPR International Workshops, SSPR 2006 and SPR 2006,

•Qian, Y., Expert, P., Rieu, T., Panzarasa, P., & Barahona, M. (2021). Quantifying the alignment of graph and features in deep learning. IEEE TNNLS

•Rong, Y., Huang, W., Xu, T., & Huang, J. (2020). Dropedge: Towards deep graph convolutional networks on node classification. In ICLR 2020.

•Rusch, T. K., Bronstein, M. M., & Mishra, S. (2023). A survey on oversmoothing in graph neural networks. arXiv preprint arXiv:2303.10993.

•Rusch, T. K., Chamberlain, B., Rowbottom, J., Mishra, S., & Bronstein, M. (2022, June). Graph-coupled oscillator networks. In International Conference on Machine Learning (pp. 18888-18909). PMLR.

Southern, J., Di Giovanni, F., Bronstein, M., & Lutzeyer, J. F. (2024). Understanding Virtual Nodes: Oversmoothing, Oversquashing, and Node Heterogeneity. arXiv preprint arXiv:2405.13526.

Shao, Z., Shi, D., Han, A., Guo, Y., Zhao, Q., & Gao, J. (2023). Unifying over-smoothing and over-squashing in graph neural networks: A physics informed approach and beyond. arXiv preprint arXiv:2309.02769.

Spielman D. (2018). Spectral Graph Theory, Lecture 10: Random Walks on Graphs. Lecture at Yale <a href="https://www.cs.yale.edu/homes/spielman/561/lect10-18.pdf">https://www.cs.yale.edu/homes/spielman/561/lect10-18.pdf</a>

•Topping, J., Di Giovanni, F., Chamberlain, B. P., Dong, X., & Bronstein, M. M. (2021). Understanding over-squashing and bottlenecks on graphs via curvature. ICLR 2022.

•Veličković, P., Cucurull, G., Casanova, A., Romero, A., Lio, P., & Bengio, Y. (2017). Graph attention networks.

 Xhonneux, L. P., Qu, M., & Tang, J. (2020, November). Continuous graph neural networks. In International conference on machine learning (pp. 10432-10441). PMLR.

Xu, K., Li, C., Tian, Y., Sonobe, T., Kawarabayashi, K. I., & Jegelka, S. (2018, July). Representation learning on graphs with jumping knowledge networks. In International conference on machine learning (pp. 5453-5462). PMLR.
Zheng, C., Zong, B., Cheng, W., Song, D., Ni, J., Yu, W., ... & Wang, W. (2020, November). Robust graph representation learning via neural sparsification. In International Conference on Machine Learning (pp. 11458-11468). PMLR.

•Zheng, Y., Luan, S., & Chen, L. (2024). What Is Missing In Homophily? Disentangling Graph Homophily For Graph Neural Networks. arXiv preprint arXiv:2406.18854.

-Zhao, L., & Akoglu, L. (2020). Pairnorm: Tackling oversmoothing in gnns. In ICLR 2020.

Zhao, J., Dong, Y., Tang, J., Ding, M., & Wang, K. (2021). Generalizing graph convolutional networks via heat kernel.
Zhou, K., Huang, X., Li, Y., Zha, D., Chen, R., & Hu, X. (2020). Towards deeper graph neural networks with

differentiable group normalization. Advances in neural information processing systems, 33, 4917-4928.

•Zhou, K., Huang, X., Zha, D., Chen, R., Li, L., Choi, S. H., & Hu, X. (2021). Dirichlet energy constrained learning for deep graph neural networks. Advances in Neural Information Processing Systems, 34, 21834-21846.

•Zhou, K., Dong, Y., Wang, K., Lee, W. S., Hooi, B., Xu, H., & Feng, J. (2021b, October). Understanding and resolving performance degradation in deep graph convolutional networks. In Proceedings of the 30th ACM International Conference on Information & Knowledge Management (pp. 2728-2737).

Zhu, J., et al. "Beyond homophily in graph neural networks: Current limitations and effective designs". in NeurIPS, 2020
Zhu, Y., Xu, W., Zhang, J., Du, Y., Zhang, J., Liu, Q., ... & Wu, S. (2021). A survey on graph structure learning: Progress and opportunities. arXiv preprint arXiv:2103.03036.

•Zhu, Y., Du, Y., Wang, Y., Xu, Y., Zhang, J., Liu, Q., & Wu, S. (2022, December). A survey on deep graph generation: Methods and applications. In Learning on Graphs Conference (pp. 47-1). PMLR.