



Università
della
Svizzera
italiana



Istituto
Dalle Molle
di studi
sull'intelligenza
artificiale

Graph Neural Networks

Operators and Architectures

Daniele Grattarola

Dissertation defence - December 10, 2021

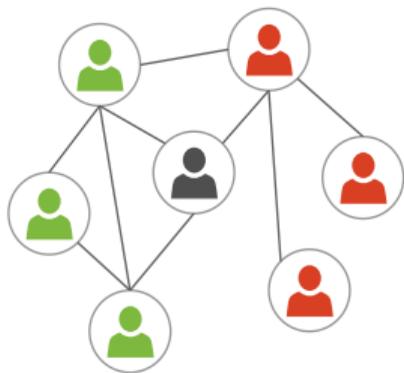
Advisor: Prof. Alippi

Co-advisor: Prof. Livi (U. Manitoba)

Internal committee: Profs. Crestani, Gambardella

External committee: Profs. Angelov (Lancaster U.), Panayiotou (U. Cyprus), Sperduti (U. Padova)

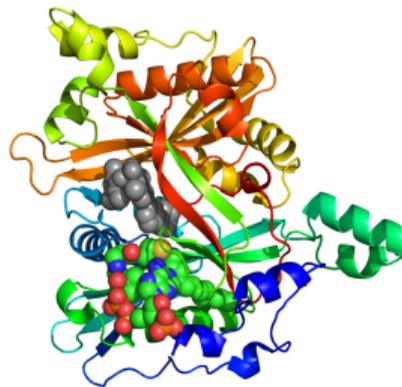
Graphs are everywhere



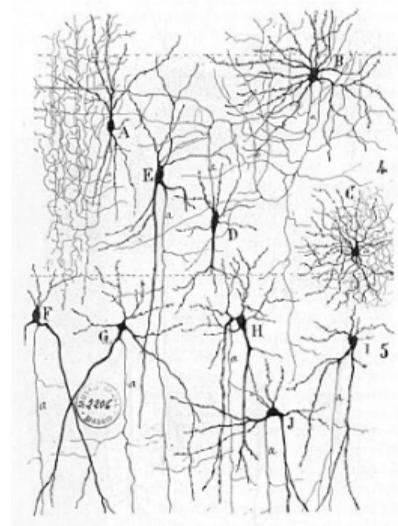
Social networks



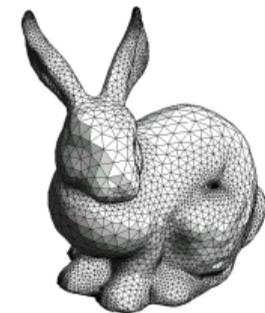
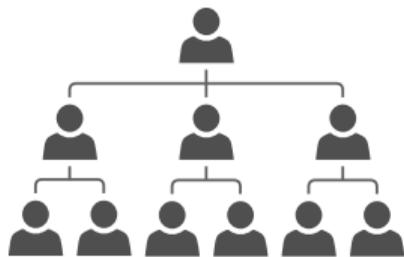
Engineering



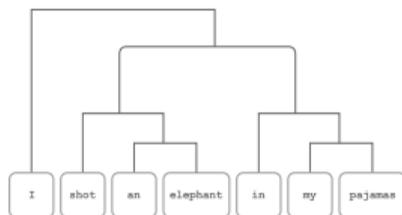
Biology



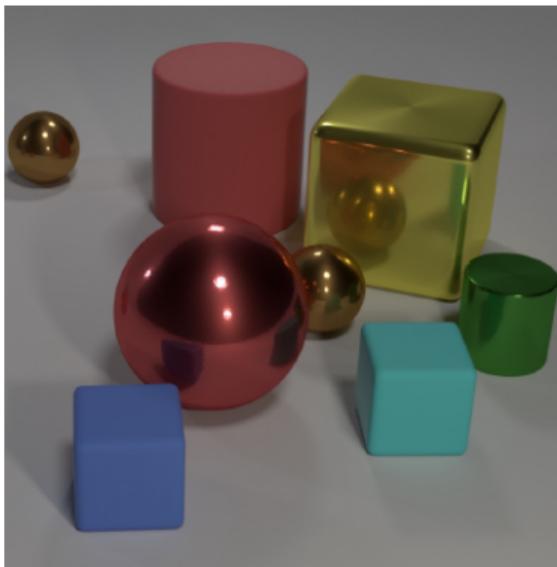
Neuroscience



Computer vision



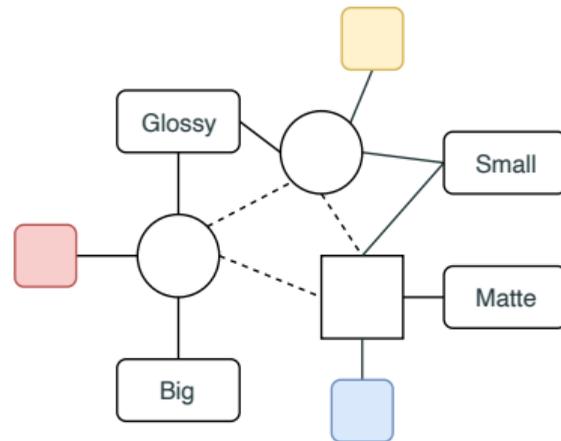
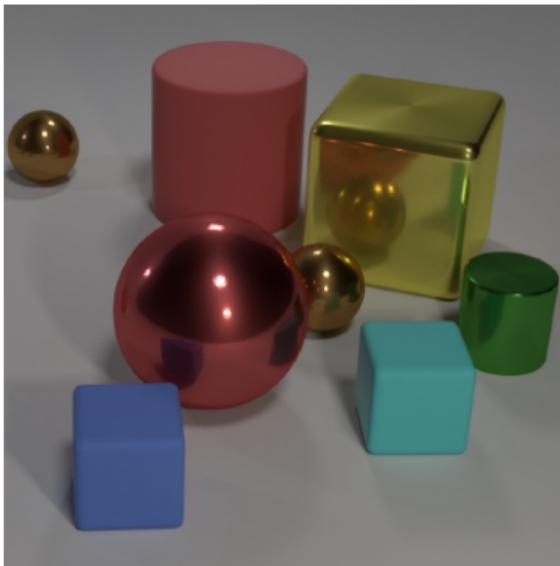
Language



Combinatorial generalisation

Combining known concepts to represent new ones.

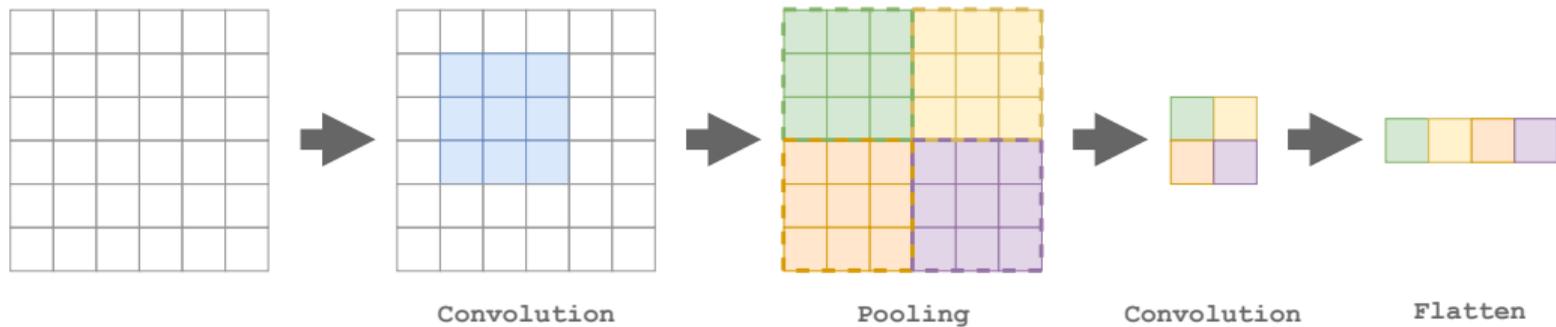
Graphs are everywhere



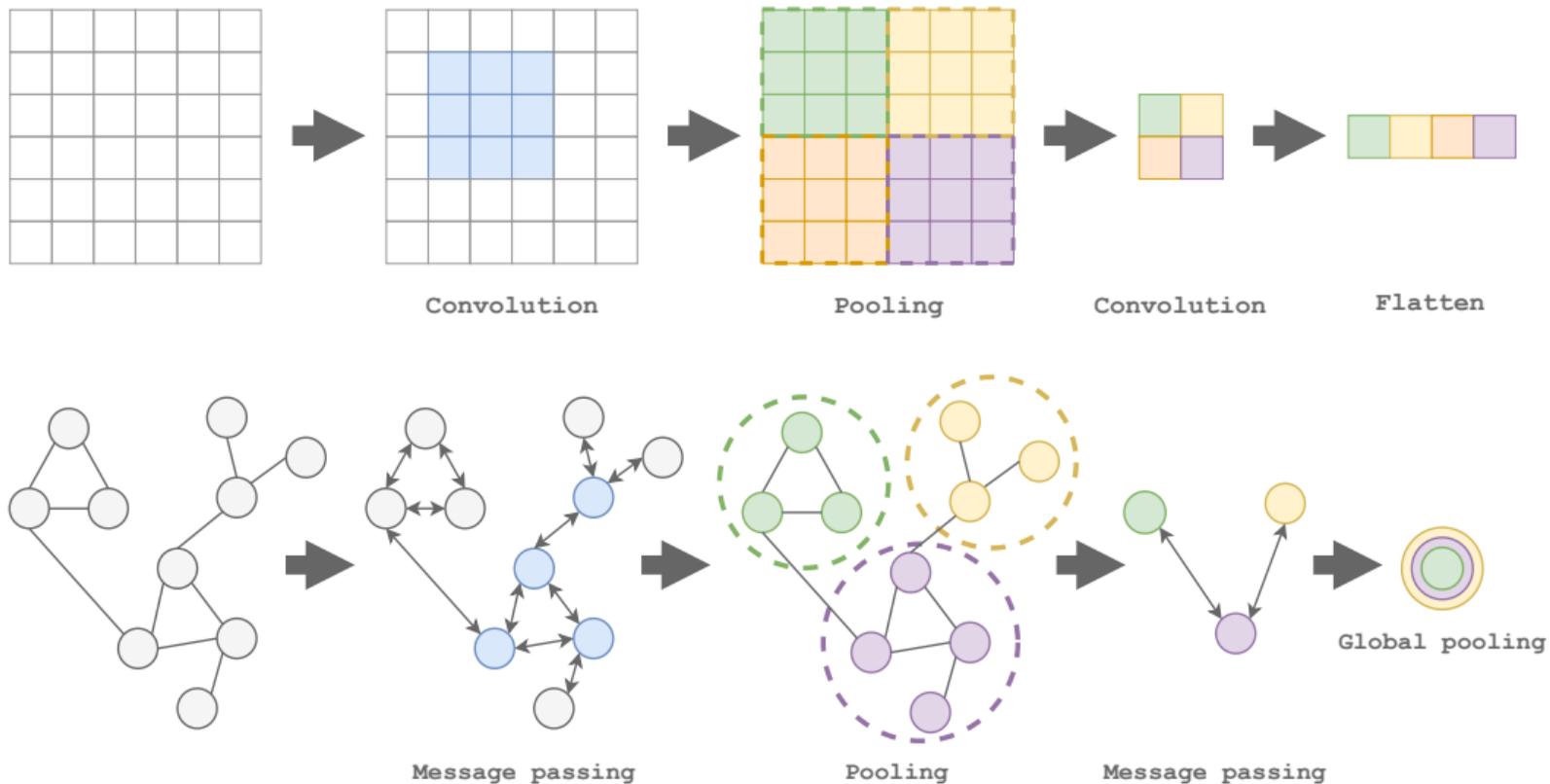
Combinatorial generalisation

Combining known concepts to represent new ones.

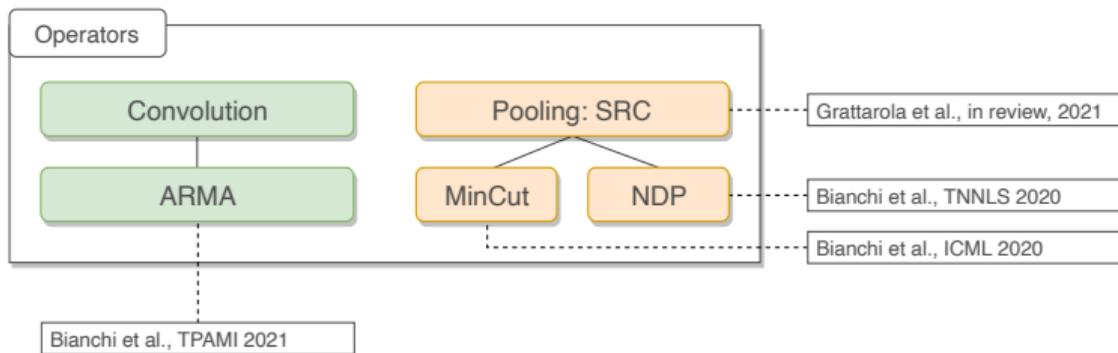
From CNNs to GNNs



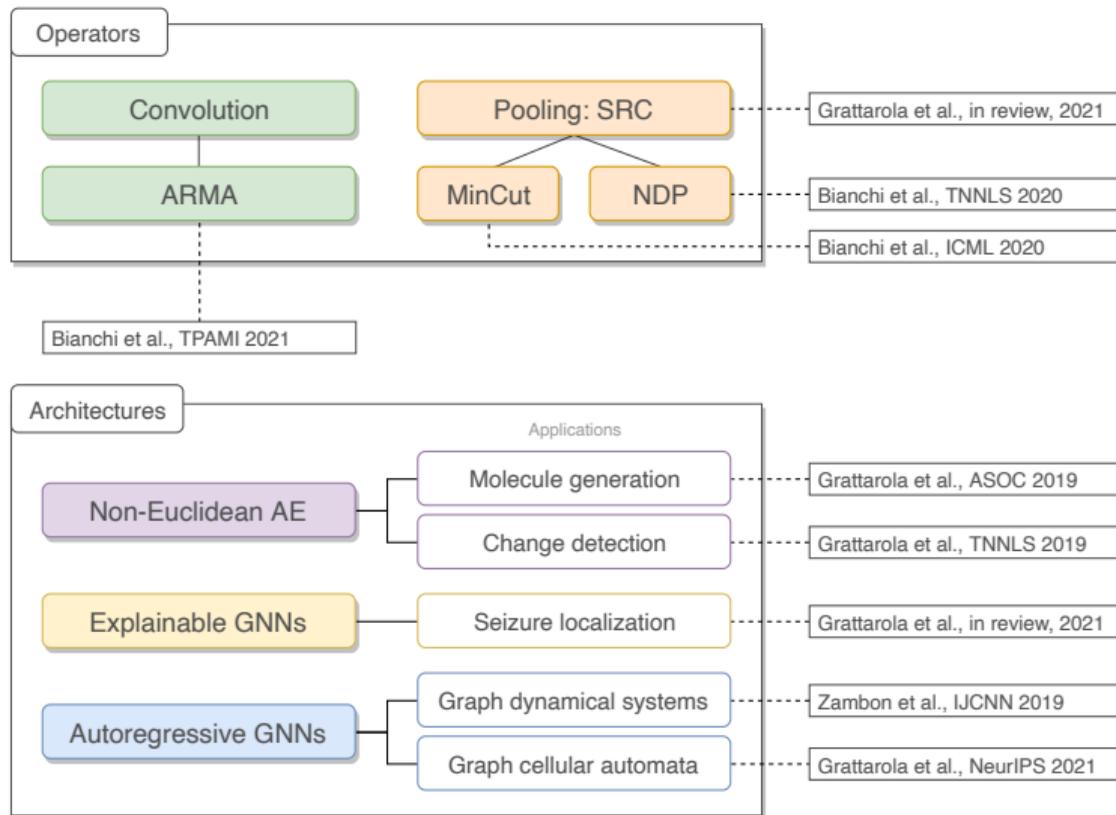
From CNNs to GNNs



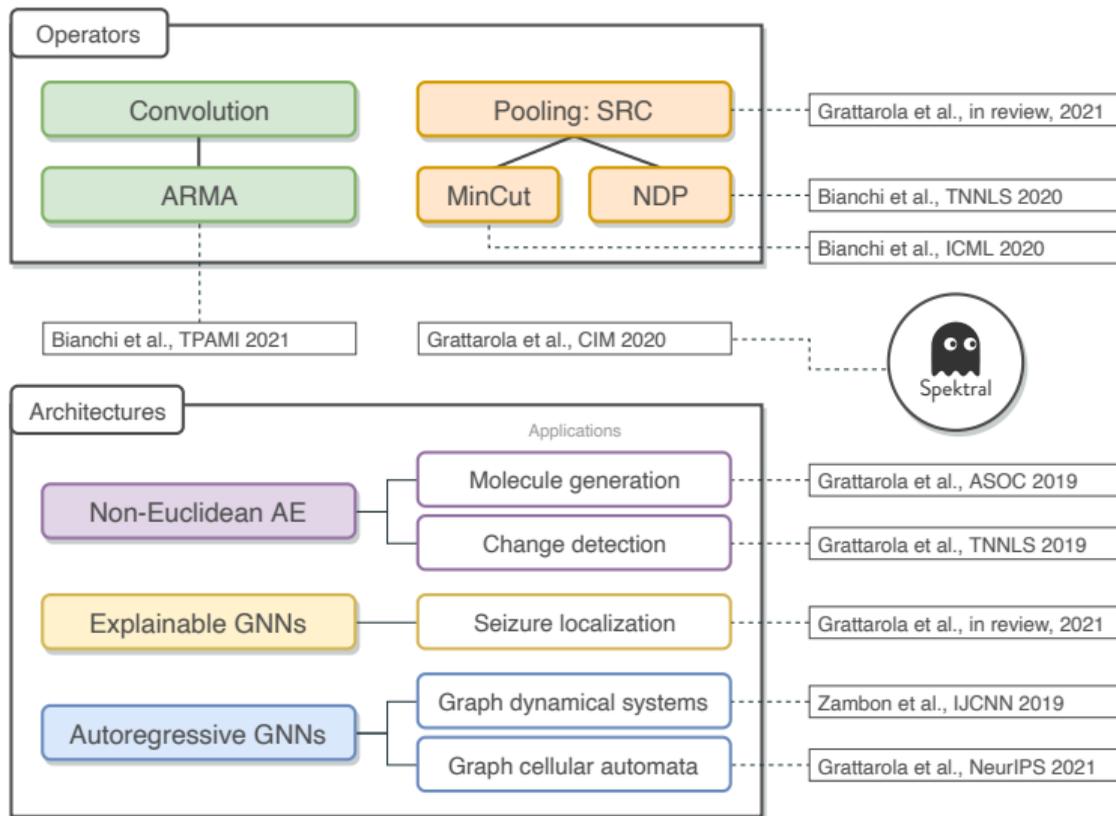
Research summary



Research summary



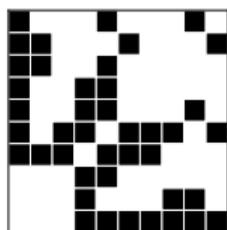
Research summary



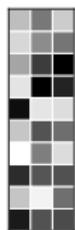
Operators

Notation

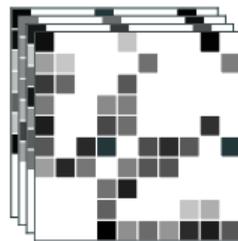
- **Graph**: nodes connected by edges
- X : node attributes
- E : edge attributes
- A : adjacency matrix
- **Structure operator** S : $s_{ij} \neq 0 \iff a_{ij} \neq 0$



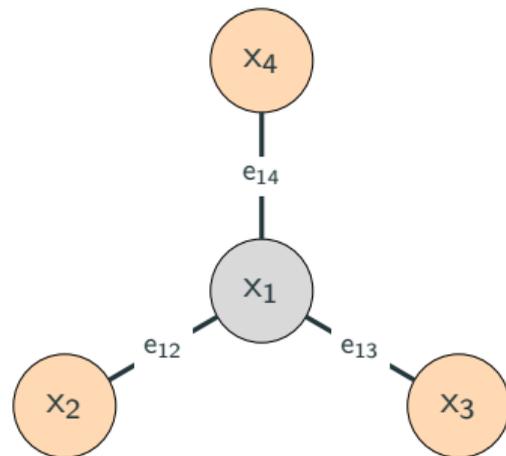
Adjacency matrix
 $N \times N$



Node attributes
 $N \times D_n$



Edge attributes
 $N \times N \times D_e$



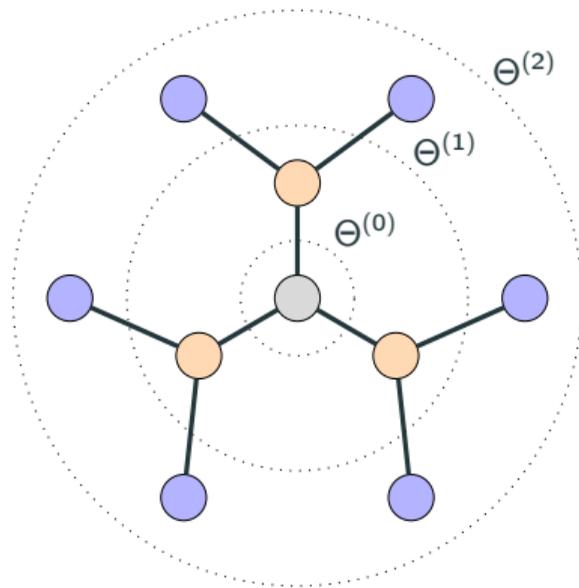
Convolutional operators

Convolutional operators - Polynomial

Typical graph convolutional networks compute a polynomial of the structure operator S [2], [3]:

$$X' = \sum_{k=0}^K S^k X \Theta^{(k)}$$

with learnable weights $\Theta^{(k)} \in \mathbb{R}^{D_n \times D'_n}$.



[2] M. Defferrard et al., "Convolutional neural networks on graphs with fast localized spectral filtering," *Advances in Neural Information Processing Systems*, 2016.

[3] T. N. Kipf et al., "Semi-Supervised Classification with Graph Convolutional Networks," *International Conference on Learning Representations (ICLR)*, 2017.

Convolution with rational filter:

$$X' = \left(I + \sum_{k=1}^K q_k S^k \right)^{-1} \left(\sum_{k=0}^{K-1} p_k S^k \right) X$$

[4] E. Isufi et al., "Autoregressive moving average graph filtering," *arXiv preprint arXiv:1602.04436*, 2016.

Convolution with rational filter:

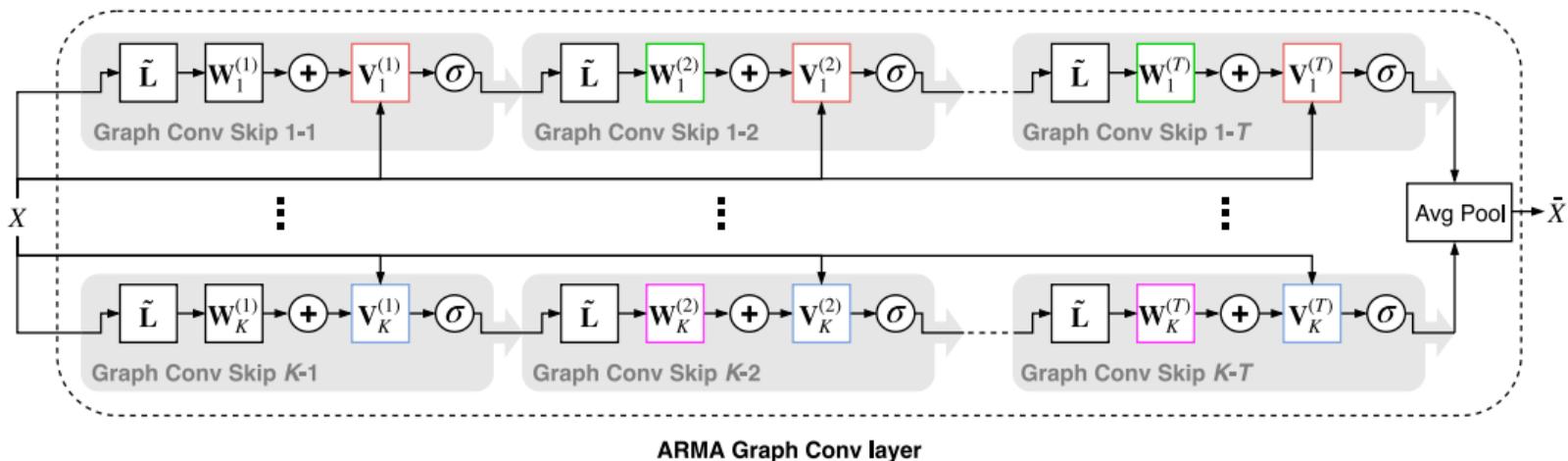
$$X' = \left(I + \sum_{k=1}^K q_k S^k \right)^{-1} \left(\sum_{k=0}^{K-1} p_k S^k \right) X$$

ARMA approximation of a rational filter [4]:

$$X^{(t+1)} = aSX^{(t)} + bX$$

[4] E. Isufi et al., "Autoregressive moving average graph filtering," *arXiv preprint arXiv:1602.04436*, 2016.

ARMA GNNs [5]



Approximate recursion with finite number of propagation steps: $X^{(t+1)} = \sigma(SX^{(t)}W + XV)$

[5] F. M. Bianchi, D. Grattarola, et al., "Graph neural networks with convolutional arma filters," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021.

Table 1: Node classification accuracy.

Method	Cora	Citeseer	Pubmed	PPI
GAT	83.1 \pm 0.6	70.9 \pm 0.6	78.5 \pm 0.3	81.3 \pm 0.1
GraphSAGE	73.7 \pm 1.8	65.9 \pm 0.9	78.5 \pm 0.6	70.0 \pm 0.0
GIN	75.1 \pm 1.7	63.1 \pm 2.0	77.1 \pm 0.7	78.1 \pm 2.6
GCN	81.5 \pm 0.4	70.1 \pm 0.7	79.0 \pm 0.5	80.8 \pm 0.1
Chebyshev	79.5 \pm 1.2	70.1 \pm 0.8	74.4 \pm 1.1	86.4 \pm 0.1
CayleyNet	81.2 \pm 1.2	67.1 \pm 2.4	75.6 \pm 3.6	84.9 \pm 1.2
ARMA	83.4 \pm 0.6	72.5 \pm 0.4	78.9 \pm 0.3	90.5 \pm 0.3

Table 2: Graph classification accuracy.

Method	Enzymes	Proteins	D&D	MUTAG	BHard
GAT	51.7 \pm 4.3	72.3 \pm 3.1	70.9 \pm 4.0	87.3 \pm 5.3	30.1 \pm 0.7
GraphSAGE	60.3 \pm 7.1	70.2 \pm 3.9	73.6 \pm 4.1	85.7 \pm 4.7	71.8 \pm 1.0
GIN	45.7 \pm 7.7	71.4 \pm 4.5	71.2 \pm 5.4	86.3 \pm 9.1	72.1 \pm 1.1
GCN	53.0 \pm 5.3	71.0 \pm 2.7	74.7 \pm 3.8	85.7 \pm 6.6	71.9 \pm 1.2
Chebyshev	57.9 \pm 2.6	72.1 \pm 3.5	73.7 \pm 3.7	82.6 \pm 5.2	71.3 \pm 1.2
CayleyNet	43.1 \pm 10.7	65.6 \pm 5.7	70.3 \pm 11.6	87.8 \pm 10.0	70.7 \pm 2.4
ARMA	60.6 \pm 7.2	73.7 \pm 3.4	77.6 \pm 2.7	91.5 \pm 4.2	74.1 \pm 0.5

Table 3: Graph signal classification accuracy.

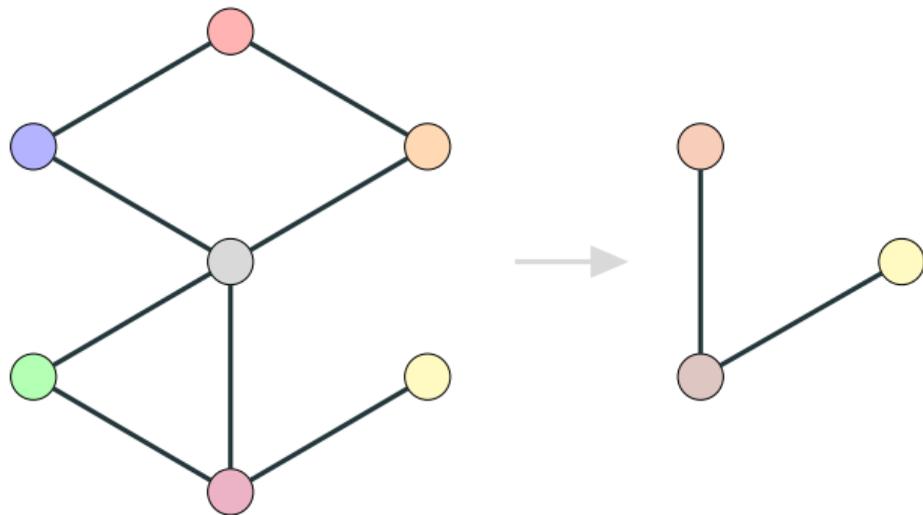
GNN layer	MNIST	20news
GCN	98.48 \pm 0.2	65.45 \pm 0.2
Chebyshev	99.14 \pm 0.1	68.24 \pm 0.2
CayleyNet	99.18 \pm 0.1	68.84 \pm 0.3
ARMA	99.20 \pm 0.1	70.02 \pm 0.1

Table 4: Graph regression mean squared error.

Property	GCN	Chebyshev	CayleyNet	ARMA
mu	0.445 \pm 0.007	0.433 \pm 0.003	0.442 \pm 0.009	0.394 \pm 0.005
alpha	0.141 \pm 0.016	0.171 \pm 0.008	0.118 \pm 0.005	0.098 \pm 0.005
HOMO	0.371 \pm 0.030	0.391 \pm 0.012	0.336 \pm 0.007	0.326 \pm 0.010
LUMO	0.584 \pm 0.051	0.528 \pm 0.005	0.679 \pm 0.148	0.508 \pm 0.011
gap	0.650 \pm 0.070	0.565 \pm 0.015	0.758 \pm 0.106	0.552 \pm 0.013
R2	0.132 \pm 0.005	0.294 \pm 0.022	0.185 \pm 0.043	0.119 \pm 0.019
ZPVE	0.349 \pm 0.022	0.358 \pm 0.001	0.555 \pm 0.174	0.338 \pm 0.001
U0_atom	0.064 \pm 0.003	0.126 \pm 0.017	1.493 \pm 1.414	0.053 \pm 0.004
Cv	0.192 \pm 0.012	0.215 \pm 0.010	0.184 \pm 0.009	0.163 \pm 0.007

Pooling operators

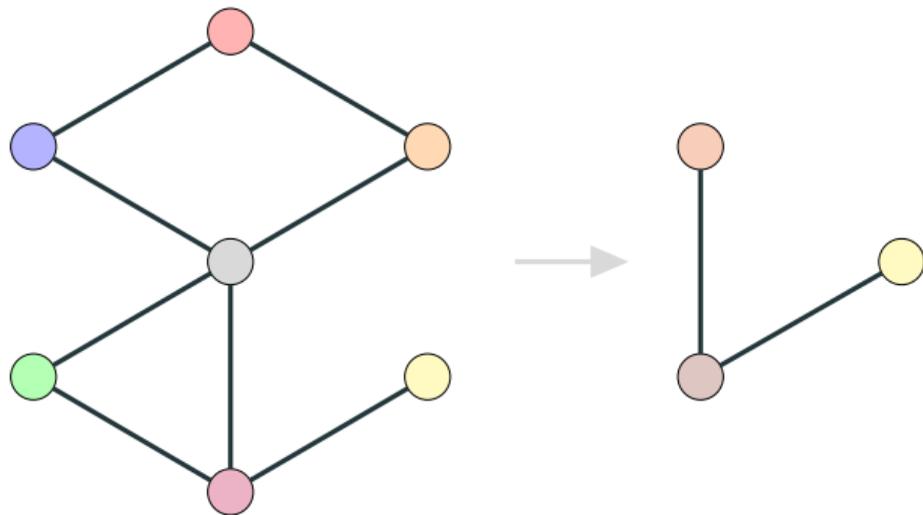
Pooling operators



Goal: reduce the size of the graph.

[6] D. Grattarola et al., "Understanding Pooling in Graph Neural Networks," *Under review at IEEE TNMLS*, 2021.

Pooling operators



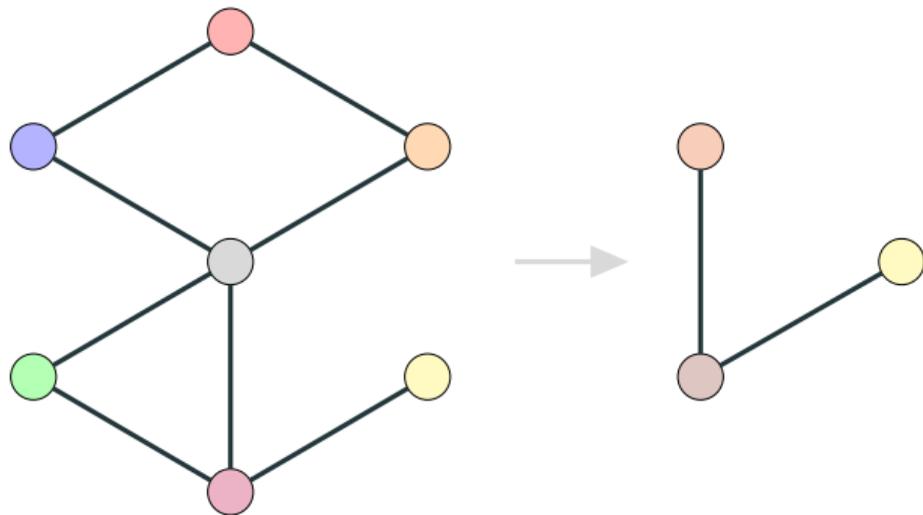
Goal: reduce the size of the graph.

Why:

- Computational cost
- Invariance to feature location
- Abstraction

[6] D. Grattarola et al., "Understanding Pooling in Graph Neural Networks," *Under review at IEEE TNMLS*, 2021.

Pooling operators



Goal: reduce the size of the graph.

Why:

- Computational cost
- Invariance to feature location
- Abstraction

Idea: 3-step process proposed in [6]:
Select, Reduce, Connect (SRC).

[6] D. Grattarola et al., "Understanding Pooling in Graph Neural Networks," *Under review at IEEE TNMLS*, 2021.

Pooling operators as SRC

Method	Select	Reduce	Connect
DiffPool [7]	$S^\top = \text{GNN}_1(A, X)$ (w/ auxiliary loss)	$X' = S \cdot \text{GNN}_2(A, X)$	$A' = SAS^\top$
MinCut [8]	$S^\top = \text{MLP}(X)$ (w/ auxiliary loss)	$X' = SX$	$A' = SAS^\top$
NMF [9]	Factorise: $A = WH \rightarrow S = H$	$X' = SX$	$A' = SAS^\top$
LaPool [10]	$\begin{cases} V = \ LX\ _d; \\ i = \{i \mid V_i > V_j, \forall j \in \mathcal{N}(i)\} \\ S^\top = \text{SparseMax} \left(\beta \frac{XX_i^\top}{\ X\ \ X_i\ } \right) \end{cases}$	$X' = SX$	$A' = SAS^\top$
Graclus [11]	$\mathcal{S}_k = \left\{ x_i, x_j \mid \arg \max_j \left(\frac{A_{ij}}{D_{ii}} + \frac{A_{ji}}{D_{jj}} \right) \right\}$	$X' = SX$	METIS [12]
NDP [13]	$i = \{i \mid \mathbf{u}_{\max, i} > 0\}$	$X' = X_i$	Kron r. [14]
Top-K [15]	$y = \frac{Xp}{\ p\ }; i = \text{top}_K(y)$	$X' = (X \odot \sigma(y))_i;$	$A' = A_{i,i}$
SAGPool [16]	$y = \text{GNN}(A, X); i = \text{top}_K(y)$	$X' = (X \odot \sigma(y))_i;$	$A' = A_{i,i}$

Selection

The **selection** stage computes K **supernodes**: $\{ \text{blue circle, red circle, orange circle} \} \{ \text{grey circle, green circle, pink circle} \} \{ \text{yellow circle} \}$

$$\text{Sel} : \mathcal{G} \mapsto \mathcal{S} = \{\mathcal{S}_1, \dots, \mathcal{S}_K\},$$

$$\mathcal{S}_k = \{(x_i, s_i) \mid s_i \in \mathbb{R}_{>0}\}.$$

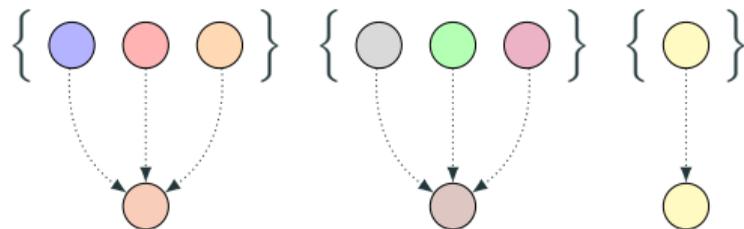
$$S = \begin{array}{|c|c|c|c|c|c|c|} \hline \text{blue} & \text{red} & \text{orange} & & & & \\ \hline & & & \text{grey} & \text{green} & \text{pink} & \\ \hline & & & & & & \text{yellow} \\ \hline \end{array} \in \mathbb{R}^{K \times N}$$

Reduction

The **reduction** stage aggregates the supernodes in a **permutation-invariant** way:

$$\text{Red} : \mathcal{G}, \mathcal{S}_k \mapsto \mathbf{x}'_k$$

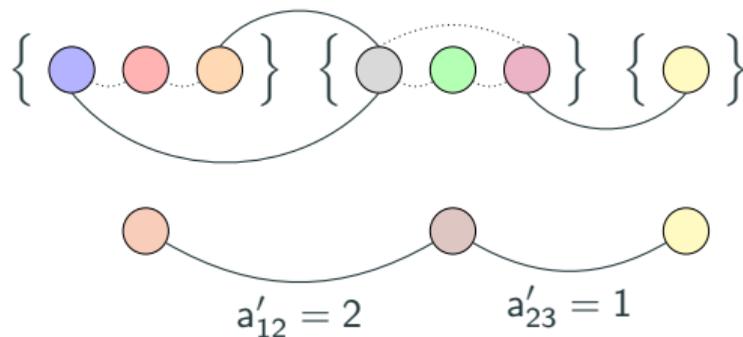
Typical approach: $X' = SX \ (\in \mathbb{R}^{K \times F})$



The **connection** function decides whether two supernodes are connected:

$$\text{Con} : \mathcal{G}, \mathcal{S}_k, \mathcal{S}_l \mapsto e'_{kl}$$

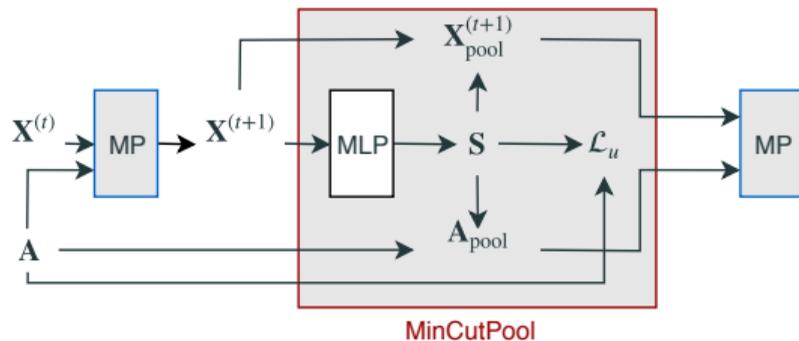
Typical approach: $A' = SAS^T$ ($\in \mathbb{R}^{K \times K}$)



MinCut pooling [8]

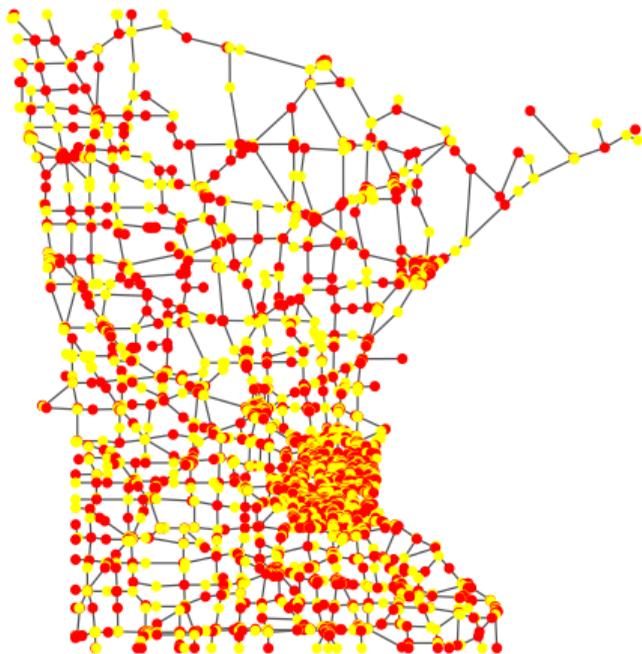
- Select: $S^T = \text{MLP}(X)$
- Reduce: $X' = SX$
- Connect: $A' = SAS^T$
- Custom loss:

$$\mathcal{L} = \underbrace{-\frac{\text{Tr}(SAS^T)}{\text{Tr}(SDS^T)}}_{\text{MinCut}} + \underbrace{\left\| \frac{SS^T}{\|SS^T\|_F} - \frac{I_K}{\sqrt{K}} \right\|_F}_{\text{Orthogonality}}$$



[8] F. M. Bianchi, D. Grattarola, et al., "Spectral Clustering with Graph Neural Networks for Graph Pooling," *International Conference on Machine Learning*, 2020.

Node decimation pooling [13]



- **Idea:** regular subsampling of the nodes using highest-frequency eigenvector u_{\max} .
- Select: $I = \{i \mid u_{\max}[i] > 0\}$
- Reduce: $X' = X_I$
- Connect: Kron reduction [14]

[14] F. Dorfler et al., "Kron Reduction of Graphs With Applications to Electrical Networks," *IEEE Transactions on Circuits and Systems I: Regular Papers*, vol. 60, no. 1, 2013.

[13] F. M. Bianchi, D. Grattarola, et al., "Hierarchical representation learning in graph neural networks with node decimation pooling," *IEEE Transactions on Neural Networks and Learning Systems*, 2020.

Table 5: Graph classification accuracy.

Dataset	WL	Dense	No-pool	Graclus	NDP	DiffPool	Top-K	SAGpool	MinCut
Bench-easy	92.6 \pm 0.0	29.3 \pm 0.3	98.5 \pm 0.3	97.5 \pm 0.5	97.9 \pm 0.5	98.6 \pm 0.4	82.4 \pm 8.9	84.2 \pm 2.3	99.0\pm0.0
Bench-hard	60.0 \pm 0.0	29.4 \pm 0.3	67.6 \pm 2.8	69.0 \pm 1.5	72.6\pm0.9	69.9 \pm 1.9	42.7 \pm 15.2	37.7 \pm 14.5	73.8\pm1.9
Mutagen.	81.7\pm1.1	68.4 \pm 0.3	78.0 \pm 1.3	74.4 \pm 1.8	77.8 \pm 2.3	77.6 \pm 2.7	71.9 \pm 3.7	72.4 \pm 2.4	79.9 \pm 2.1
Proteins	71.2 \pm 2.6	68.7 \pm 3.3	72.6 \pm 4.8	68.6 \pm 4.6	73.3 \pm 3.7	72.7 \pm 3.8	69.6 \pm 3.5	70.5 \pm 2.6	76.5\pm2.6
DD	78.6 \pm 2.7	70.6 \pm 5.2	76.8 \pm 1.5	70.5 \pm 4.8	72.0 \pm 3.1	79.3\pm2.4	69.4 \pm 7.8	71.5 \pm 4.5	80.8\pm2.3
COLLAB	74.8 \pm 1.3	79.3 \pm 1.6	82.1\pm1.8	77.1 \pm 2.1	79.1 \pm 1.5	81.8 \pm 1.4	79.3 \pm 1.8	79.2 \pm 2.0	83.4\pm1.7
Reddit-B	68.2 \pm 1.7	48.5 \pm 2.6	80.3 \pm 2.6	79.2 \pm 0.4	84.3 \pm 2.4	86.8 \pm 2.1	74.7 \pm 4.5	73.9 \pm 5.1	91.4\pm1.5

MinCut achieves state of the art on many graph classification benchmarks (results from [8]).

[8] F. M. Bianchi, D. Grattarola, et al., "Spectral Clustering with Graph Neural Networks for Graph Pooling," *International Conference on Machine Learning*, 2020.

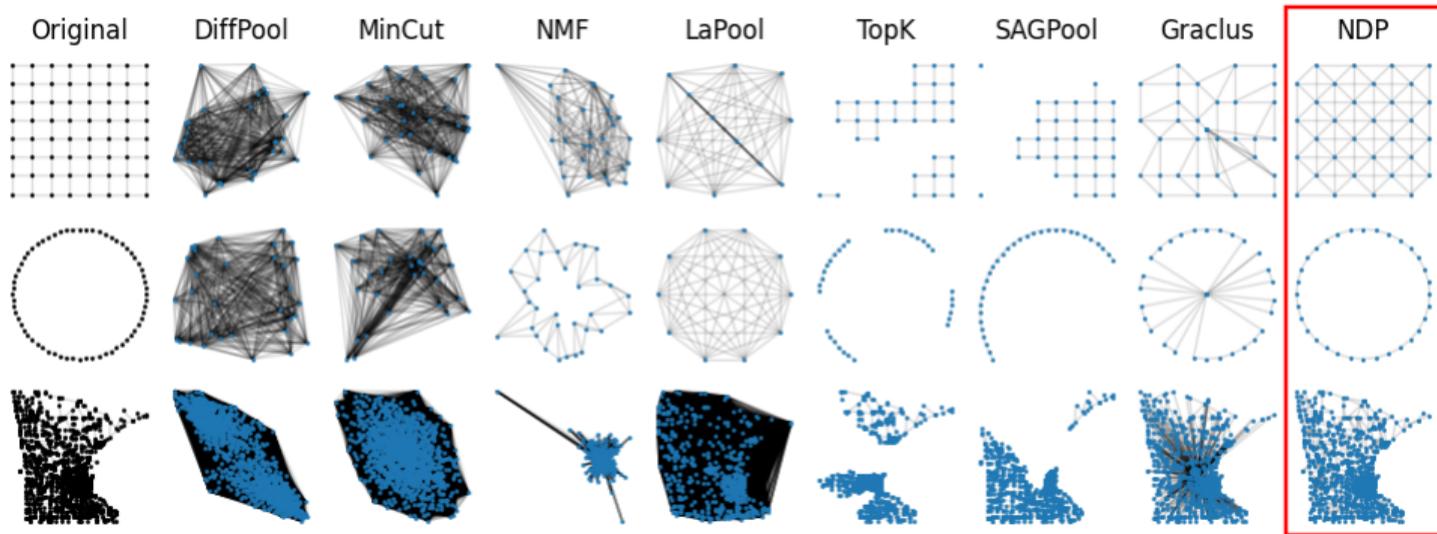
Table 6: Reconstruction MSE (scale of 10^{-3}).

	DiffPool	MinCut	NMF	LaPool	TopK	SAGPool	Graclus	NDP
Grid2d	0.010 ± 0.005	0.002 ± 0.002	0.000 ± 0.000	0.002 ± 0.001	18.86 ± 3.923	16.61 ± 3.270	0.109 ± 0.000	0.000 ± 0.000
Ring	0.018 ± 0.003	0.001 ± 0.000	0.000 ± 0.000	0.052 ± 0.046	132.2 ± 4.133	148.5 ± 30.10	0.600 ± 0.000	0.000 ± 0.000
Bunny	3.901 ± 0.275	0.208 ± 0.034	0.339 ± 0.055	0.610 ± 0.103	15.32 ± 3.557	16.10 ± 1.722	0.332 ± 0.043	0.373 ± 0.070
Airplane	0.094 ± 0.022	0.005 ± 0.002	0.020 ± 0.000	0.002 ± 0.000	0.096 ± 0.028	0.268 ± 0.081	0.009 ± 0.000	0.012 ± 0.000
Car	0.143 ± 0.127	0.535 ± 0.200	0.016 ± 0.001	OR	0.229 ± 0.023	0.204 ± 0.029	0.102 ± 0.000	0.009 ± 0.000
Guitar	0.101 ± 0.025	0.313 ± 0.000	0.007 ± 0.000	OR	0.056 ± 0.051	0.060 ± 0.044	0.010 ± 0.000	0.005 ± 0.000
Person	0.077 ± 0.041	0.301 ± 0.000	0.001 ± 0.000	OR	0.055 ± 0.012	0.062 ± 0.033	0.001 ± 0.000	0.001 ± 0.000

NDP works better for geometric data, where regular subsampling is desirable (results from [6]).

[6] D. Grattarola et al., "Understanding Pooling in Graph Neural Networks," *Under review at IEEE TNMLS*, 2021.

Pooling operators - Results



NDP works better for geometric data, where regular subsampling is desirable (results from [6]).

[6] D. Grattarola et al., "Understanding Pooling in Graph Neural Networks," *Under review at IEEE TNMLS*, 2021.

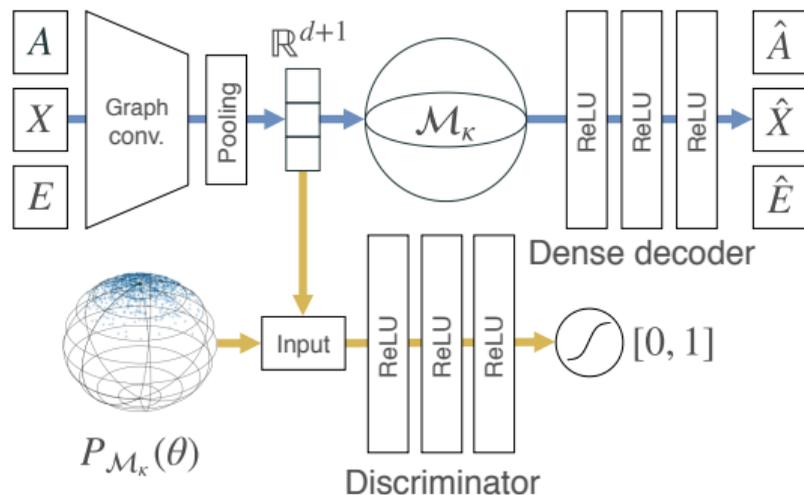
Architectures

- Non-Euclidean Autoencoders
- Explainable GNNs
- Auto-regressive GNNs

Non-Euclidean Autoencoders

Non-Euclidean Autoencoders [17]

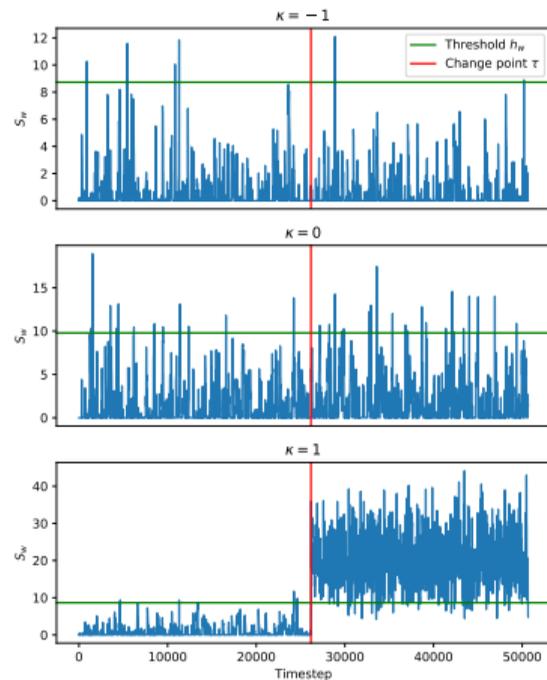
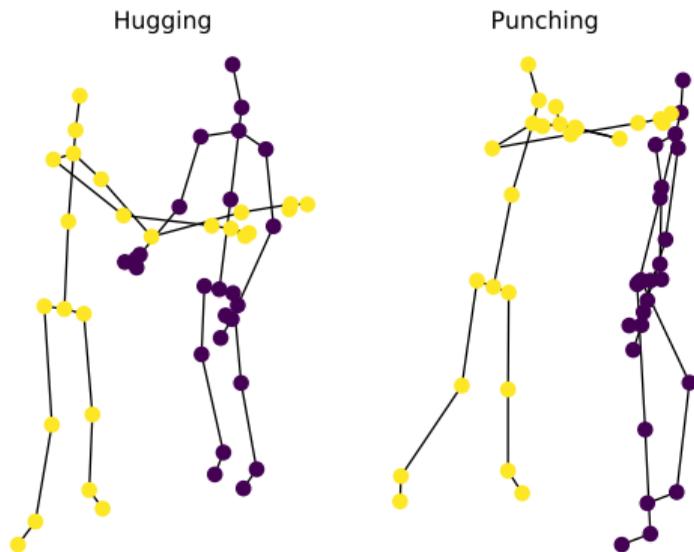
Goal: represent graphs on a non-Euclidean manifold with constant curvature (CCM).
How: adversarial AE with non-Euclidean prior.



$$\min_{f_{\text{enc}}} \max_{f_{\text{dis}}} \mathbb{E}_{z \sim p(z)} [\log f_{\text{dis}}(z)] + \mathbb{E}_{x \sim p_{\text{data}}(x)} [\log(1 - (f_{\text{dis}} \circ f_{\text{enc}})(x))].$$

[17] D. Grattarola et al., "Adversarial autoencoders with constant-curvature latent manifolds," *Applied Soft Computing*, 2019.

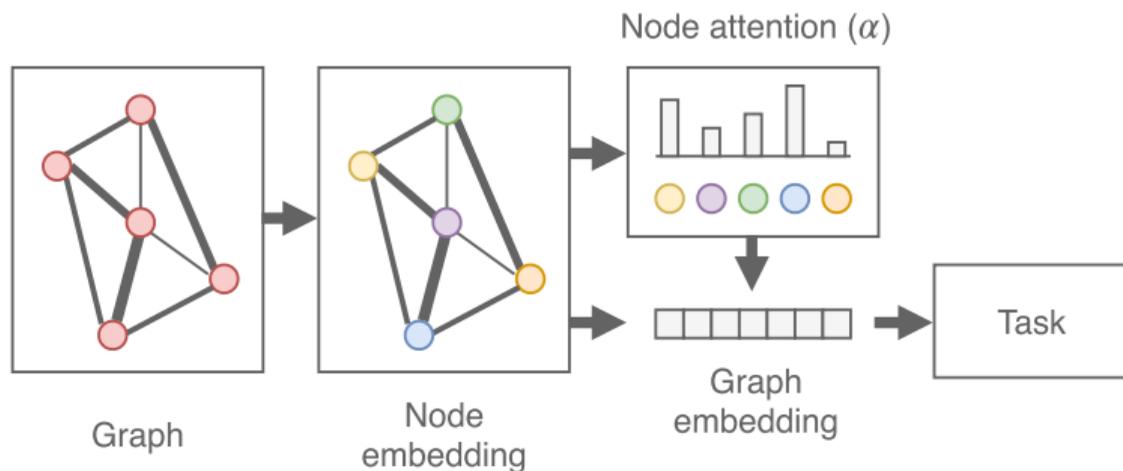
Change detection on CCMs [18]



[18] D. Grattarola et al., "Change Detection in Graph Streams by Learning Graph Embeddings on Constant-Curvature Manifolds," *IEEE Transactions on Neural Networks and Learning Systems*, 2019.

Explainable GNNs

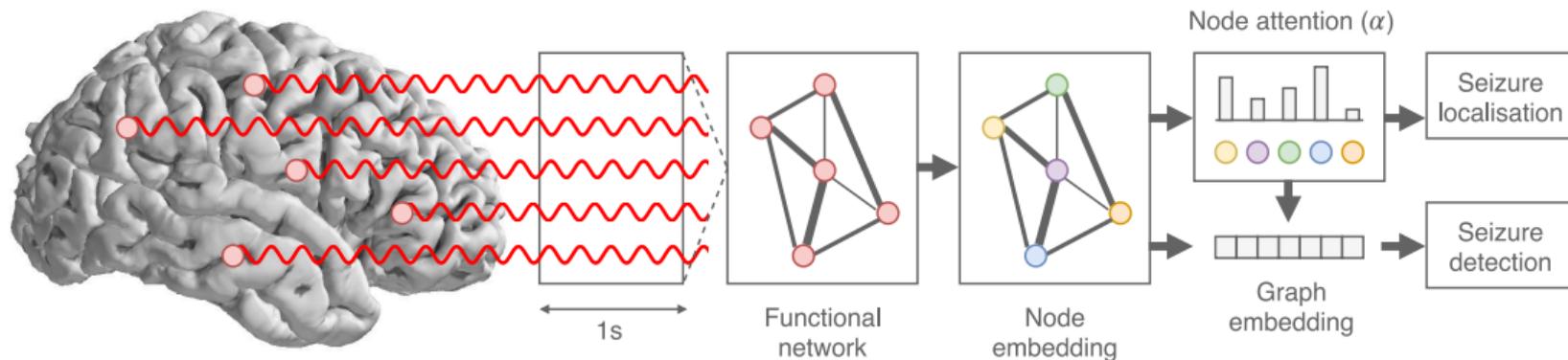
Explainable GNNs



Make GNN explainable by introducing attention in the readout operation:

$$z = \text{Attn-RO}(h) = \sum_{j=1}^N \alpha_j h_j, \text{ where } \alpha_j = \frac{\exp(h_j \cdot a)}{\sum_{k=1}^N \exp(h_k \cdot a)},$$

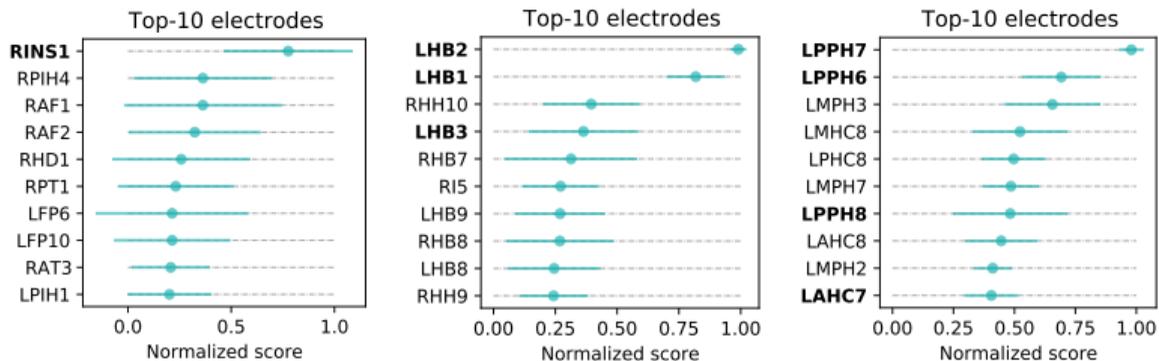
Explainable GNNs for seizure localization [19]



Idea: use attention scores to identify the brain areas where seizures originate.

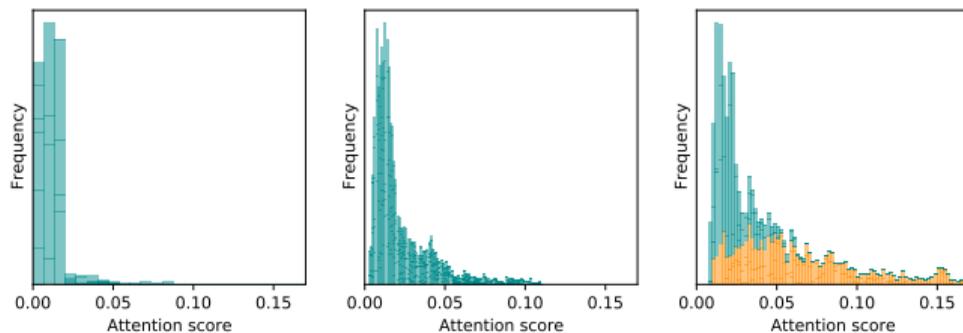
[19] D. Grattarola et al., "Unsupervised seizure localisation with attention-based graph neural networks," *Under review at IEEE TBME*, 2021.

Seizure localization - Results



Known seizure onset zone: GNN shows strong correlation with clinical diagnosis.

Seizure localization - Results



(a) Unknown SOZ

(b) Unknown SOZ

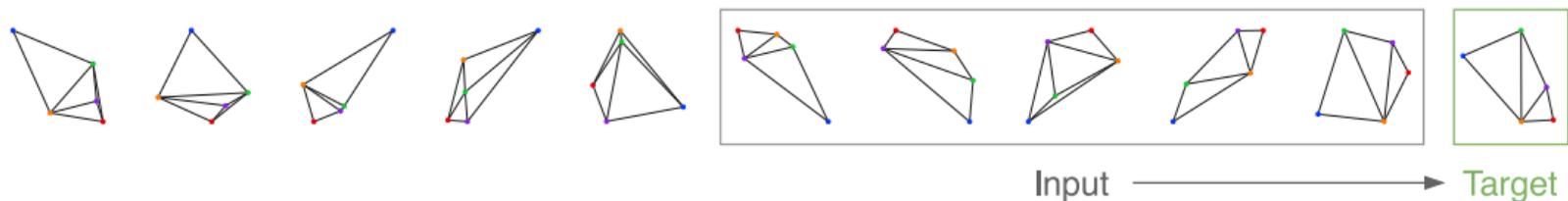
(c) Known SOZ

Unknown seizure onset zone: GNN also shows uncertainty.

Auto-regressive GNNs

Auto-regressive GNNs [20]

Goal: predict the next observation of a graph-valued dynamical system.



[20] D. Zambon *et al.*, "Autoregressive Models for Sequences of Graphs," *International Joint Conference on Neural Networks*, 2019.

Traditional AR model

$$f : \mathbb{R}^{p \times d} \rightarrow \mathbb{R}^d$$

$$x_{t+1} = f(x_t, x_{t-1}, \dots, x_{t-p+1}) + \epsilon$$

$$\mathbb{E}[\epsilon_i] = 0$$

$$\text{Var}[\epsilon_i] = \sigma^2 < \infty$$

Graph AR model

Traditional AR model

$$f : \mathbb{R}^{p \times d} \rightarrow \mathbb{R}^d$$

$$x_{t+1} = f(x_t, x_{t-1}, \dots, x_{t-p+1}) + \epsilon$$

$$\mathbb{E}[\epsilon_i] = 0$$

$$\text{Var}[\epsilon_i] = \sigma^2 < \infty$$

Graph AR model

$$\phi : \mathfrak{G}^p \rightarrow \mathfrak{G}$$

Traditional AR model

$$f : \mathbb{R}^{p \times d} \rightarrow \mathbb{R}^d$$

$$x_{t+1} = f(x_t, x_{t-1}, \dots, x_{t-p+1}) + \epsilon$$

$$\mathbb{E}[\epsilon_i] = 0$$

$$\text{Var}[\epsilon_i] = \sigma^2 < \infty$$

Graph AR model

$$\phi : \mathfrak{G}^p \rightarrow \mathfrak{G}$$

$$\mathcal{G}_{t+1} = H(\phi(\mathcal{G}_t, \dots, \mathcal{G}_{t-p+1}), \eta)$$

Traditional AR model

$$f : \mathbb{R}^{p \times d} \rightarrow \mathbb{R}^d$$

$$x_{t+1} = f(x_t, x_{t-1}, \dots, x_{t-p+1}) + \epsilon$$

$$\mathbb{E}[\epsilon_i] = 0$$

$$\text{Var}[\epsilon_i] = \sigma^2 < \infty$$

Graph AR model

$$\phi : \mathfrak{G}^p \rightarrow \mathfrak{G}$$

$$\mathcal{G}_{t+1} = H(\phi(\mathcal{G}_t, \dots, \mathcal{G}_{t-p+1}), \eta)$$

$$\phi(\mathcal{G}_t, \dots, \mathcal{G}_{t-p+1}) \in \mathbb{E}_\eta^f[H(\phi(\mathcal{G}_t, \dots, \mathcal{G}_{t-p+1}), \eta)]$$

Where:

$$\mathbb{E}^f[\mathcal{G}] = \arg \min_{\mathcal{G}' \in \mathfrak{G}} \int_{\mathfrak{G}} d(\mathcal{G}, \mathcal{G}')^2 dQ(\mathcal{G})$$

Traditional AR model

$$f : \mathbb{R}^{p \times d} \rightarrow \mathbb{R}^d$$

$$x_{t+1} = f(x_t, x_{t-1}, \dots, x_{t-p+1}) + \epsilon$$

$$\mathbb{E}[\epsilon_i] = 0$$

$$\text{Var}[\epsilon_i] = \sigma^2 < \infty$$

Graph AR model

$$\phi : \mathfrak{G}^p \rightarrow \mathfrak{G}$$

$$\mathcal{G}_{t+1} = H(\phi(\mathcal{G}_t, \dots, \mathcal{G}_{t-p+1}), \eta)$$

$$\phi(\mathcal{G}_t, \dots, \mathcal{G}_{t-p+1}) \in \mathbb{E}_\eta^f[H(\phi(\mathcal{G}_t, \dots, \mathcal{G}_{t-p+1}), \eta)]$$

$$\text{Var}^f[\eta] < \infty$$

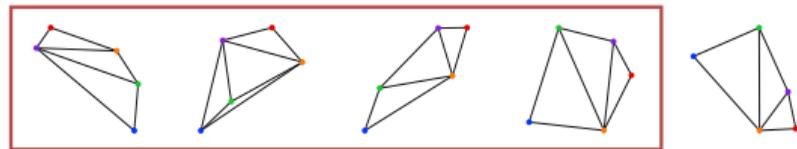
Where:

$$\mathbb{E}^f[\mathcal{G}] = \arg \min_{\mathcal{G}' \in \mathfrak{G}} \int_{\mathfrak{G}} d(\mathcal{G}, \mathcal{G}')^2 dQ(\mathcal{G}) \quad \text{Var}^f[\mathcal{G}] := \min_{\mathcal{G}' \in \mathfrak{G}} \int_{\mathfrak{G}} d(\mathcal{G}, \mathcal{G}')^2 dQ(\mathcal{G})$$

Auto-regressive GNN

Input

Window of observations up to time t .



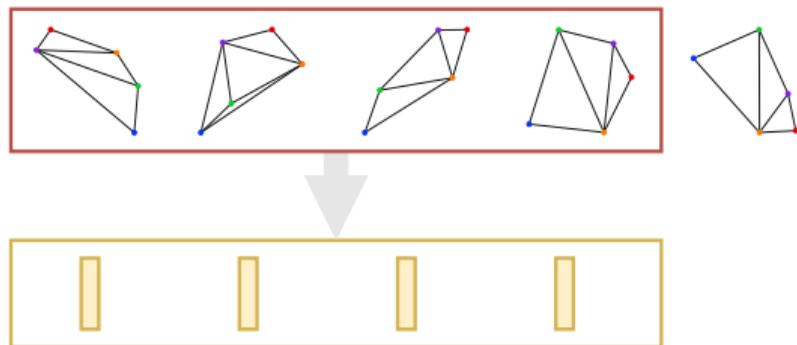
Auto-regressive GNN

Input

Window of observations up to time t .

Graph embedding

GNN block applied in parallel to each graph.



Auto-regressive GNN

Input

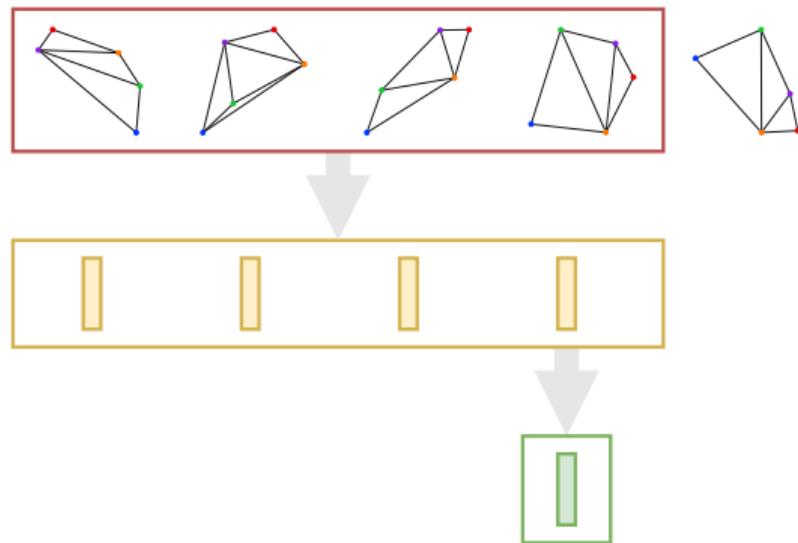
Window of observations up to time t .

Graph embedding

GNN block applied in parallel to each graph.

Predictor

RNN predicts the graph embedding at $t + 1$.



Auto-regressive GNN

Input

Window of observations up to time t .

Graph embedding

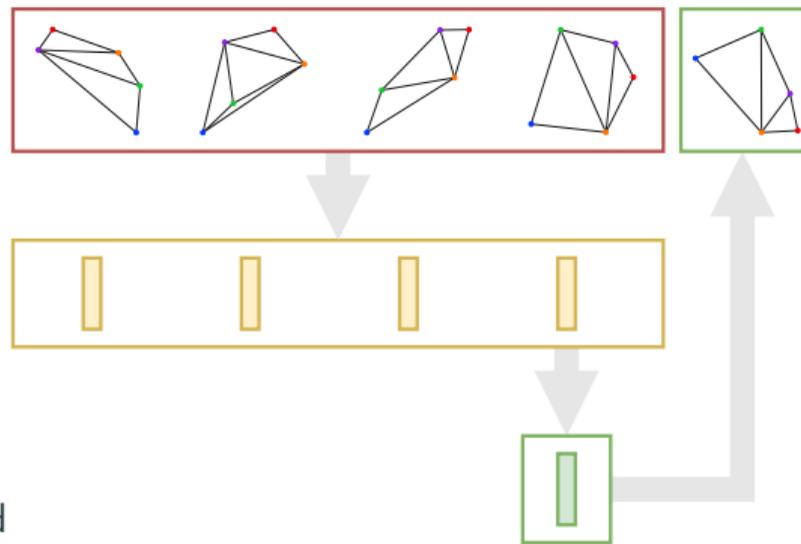
GNN block applied in parallel to each graph.

Predictor

RNN predicts the graph embedding at $t + 1$.

Graph decoder

Maps the predicted embedding to the predicted graph.



Auto-regressive GNNs - Results

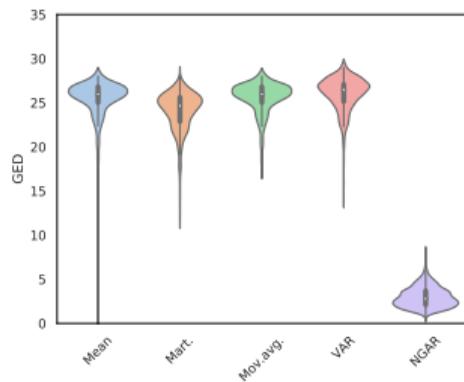
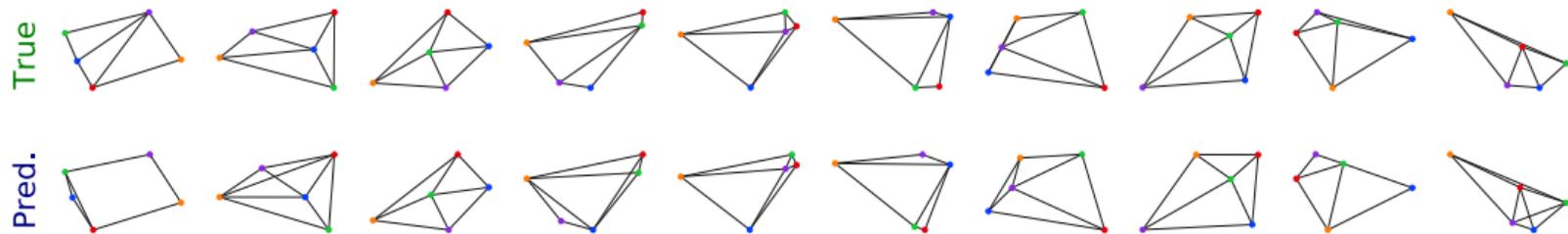


Figure 4: Distance b/w predicted and true graphs.

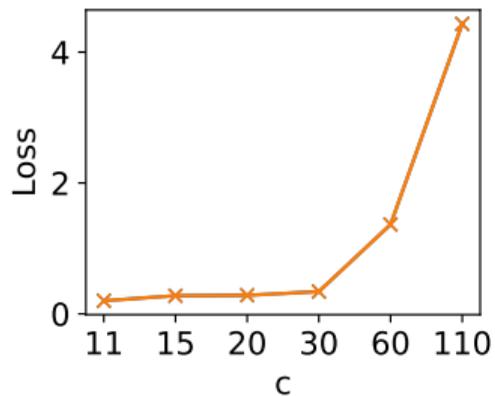


Figure 5: Loss v. complexity of the problem.

Spektral



Python library for GNNs:

- TensorFlow/Keras
- 25+ layers for convolution and pooling
- Easy to use, flexible, fast
- Almost 2000 ★ on github.com

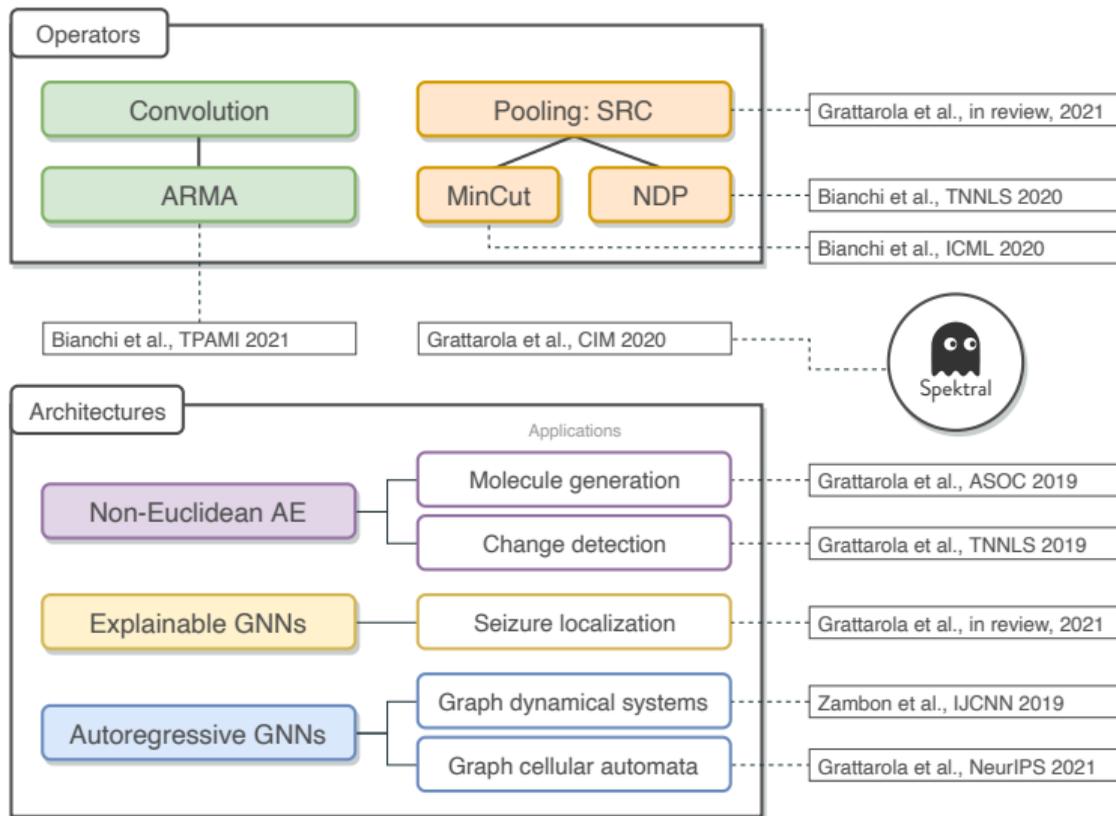
Website: graphneural.network

Featured on IEEE Computational Intelligence Magazine [21].

[21] D. Grattarola et al., "Graph neural networks in Tensorflow and Keras with Spektral," *IEEE Computational Intelligence Magazine*, 2021.

Conclusions

Summary



- [1] P. W. Battaglia, J. B. Hamrick, V. Bapst, *et al.*, “Relational inductive biases, deep learning, and graph networks,” *arXiv preprint arXiv:1806.01261*, 2018.
- [2] M. Defferrard, X. Bresson, and P. Vandergheynst, “Convolutional neural networks on graphs with fast localized spectral filtering,” *Advances in Neural Information Processing Systems*, pp. 3844–3852, 2016.
- [3] T. N. Kipf and M. Welling, “Semi-supervised classification with graph convolutional networks,” *International Conference on Learning Representations (ICLR)*, 2017.
- [4] E. Isufi, A. Loukas, A. Simonetto, and G. Leus, “Autoregressive moving average graph filtering,” *arXiv preprint arXiv:1602.04436*, 2016.
- [5] F. M. Bianchi, **D. Grattarola**, L. Livi, and C. Alippi, “Graph neural networks with convolutional arma filters,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021.

- [6] D. Grattarola, D. Zambon, F. M. Bianchi, and C. Alippi, “Understanding pooling in graph neural networks,” *Under review at IEEE TNNLS*, 2021.
- [7] R. Ying, J. You, C. Morris, X. Ren, W. L. Hamilton, and J. Leskovec, “Hierarchical graph representation learning with differentiable pooling,” *arXiv preprint arXiv:1806.08804*, 2018.
- [8] F. M. Bianchi, D. Grattarola, and C. Alippi, “Spectral clustering with graph neural networks for graph pooling,” *International Conference on Machine Learning*, 2020.
- [9] D. Bacciu and L. Di Sotto, “A non-negative factorization approach to node pooling in graph convolutional neural networks,” in *Proceedings of the 18th International Conference of the Italian Association for Artificial Intelligence, AIIA*, 2019.
- [10] E. Noutahi, D. Beani, J. Horwood, and P. Tossou, “Towards interpretable sparse graph representation learning with laplacian pooling,” *arXiv preprint arXiv:1905.11577*, 2019.

- [11] I. S. Dhillon, Y. Guan, and B. Kulis, “Weighted graph cuts without eigenvectors a multilevel approach,” *IEEE transactions on pattern analysis and machine intelligence*, vol. 29, no. 11, pp. 1944–1957, 2007.
- [12] G. Karypis, “Metis: Unstructured graph partitioning and sparse matrix ordering system,” *Technical report*, 1997.
- [13] F. M. Bianchi, D. Grattarola, L. Livi, and C. Alippi, “Hierarchical representation learning in graph neural networks with node decimation pooling,” *IEEE Transactions on Neural Networks and Learning Systems*, 2020.
- [14] F. Dorfler and F. Bullo, “Kron reduction of graphs with applications to electrical networks,” *IEEE Transactions on Circuits and Systems I: Regular Papers*, vol. 60, no. 1, pp. 150–163, 2013, issn: 1549-8328. doi: 10.1109/TCSI.2012.2215780.
- [15] S. J. Hongyang Gao, “Graph u-net,” *Submitted to ICLR*, 2019.

- [16] J. Lee, I. Lee, and J. Kang, “Self-attention graph pooling,” *arXiv preprint arXiv:1904.08082*, 2019.
- [17] D. Grattarola, L. Livi, and C. Alippi, “Adversarial autoencoders with constant-curvature latent manifolds,” *Applied Soft Computing*, 2019.
- [18] D. Grattarola, D. Zambon, C. Alippi, and L. Livi, “Change detection in graph streams by learning graph embeddings on constant-curvature manifolds,” *IEEE Transactions on Neural Networks and Learning Systems*, 2019.
- [19] D. Grattarola, L. Livi, C. Alippi, R. Wennberg, and T. Valiante, “Unsupervised seizure localisation with attention-based graph neural networks,” *Under review at IEEE TBME*, 2021.
- [20] D. Zambon, D. Grattarola, L. Livi, and C. Alippi, “Autoregressive models for sequences of graphs,” *International Joint Conference on Neural Networks*, 2019.

- [21] D. Grattarola and C. Alippi, “Graph neural networks in tensorflow and keras with spektral,” *IEEE Computational Intelligence Magazine*, 2021.
- [22] A. Sperduti and A. Starita, “Supervised neural networks for the classification of structures,” *IEEE Transactions on Neural Networks*, vol. 8, no. 3, pp. 714–735, 1997.
- [23] M. Gori, G. Monfardini, and F. Scarselli, “A new model for learning in graph domains,” in *Proceedings. 2005 IEEE International Joint Conference on Neural Networks, 2005.*, IEEE, vol. 2, 2005, pp. 729–734.
- [24] F. Scarselli, M. Gori, A. C. Tsoi, M. Hagenbuchner, and G. Monfardini, “The graph neural network model,” *IEEE Transactions on Neural Networks*, vol. 20, no. 1, pp. 61–80, 2009.
- [25] J. Bruna, W. Zaremba, A. Szlam, and Y. LeCun, “Spectral networks and locally connected networks on graphs,” *arXiv preprint arXiv:1312.6203*, 2013.

- [26] M. M. Bronstein, J. Bruna, Y. LeCun, A. Szlam, and P. Vandergheynst, “Geometric deep learning: Going beyond euclidean data,” *IEEE Signal Processing Magazine*, vol. 34, no. 4, pp. 18–42, 2017.
- [27] J. Gilmer, S. S. Schoenholz, P. F. Riley, O. Vinyals, and G. E. Dahl, “Neural message passing for quantum chemistry,” *arXiv preprint arXiv:1704.01212*, 2017.
- [28] M. Simonovsky and N. Komodakis, “Dynamic edgeconditioned filters in convolutional neural networks on graphs,” in *Proc. CVPR*, 2017.
- [29] J. Shi and J. Malik, “Normalized cuts and image segmentation,” *Departmental Papers (CIS)*, p. 107, 2000.
- [30] U. Von Luxburg, “A tutorial on spectral clustering,” *Statistics and computing*, vol. 17, no. 4, pp. 395–416, 2007.

- [31] L. Palagi, V. Piccialli, F. Rendl, G. Rinaldi, and A. Wiegele, “Computational approaches to max-cut,” in *Handbook on semidefinite, conic and polynomial optimization*, Springer, 2012, pp. 821–847.
- [32] C. Bodnar, C. Cangea, and P. Liò, “Deep graph mapper: Seeing graphs through the neural lens,” *arXiv preprint arXiv:2002.03864*, 2020.
- [33] G. Singh, F. Mémoli, and G. E. Carlsson, “Topological methods for the analysis of high dimensional data sets and 3d object recognition.,” in *SPBG*, 2007, pp. 91–100.

Backup

Convolutional operators

- **1997**: first GNN (works only for DAGs) [22]
- **2005**: first use of the term "GNN" [23]
- **2009**: improved version of 2005 paper [24]
- **2013**: First GCN [25]
- **2016**: Geometric Deep Learning [26]

[22] A. Sperduti et al., "Supervised neural networks for the classification of structures," *IEEE Transactions on Neural Networks*, vol. 8, no. 3, 1997.

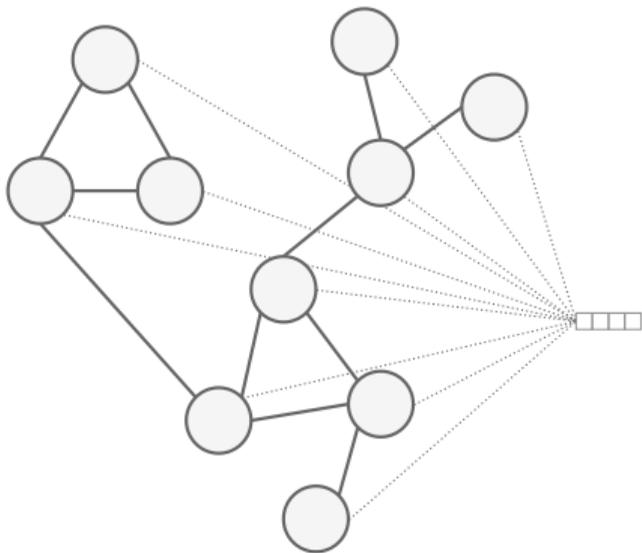
[23] M. Gori et al., "A new model for learning in graph domains," vol. 2, 2005.

[24] F. Scarselli et al., "The graph neural network model," *IEEE Transactions on Neural Networks*, vol. 20, no. 1, 2009.

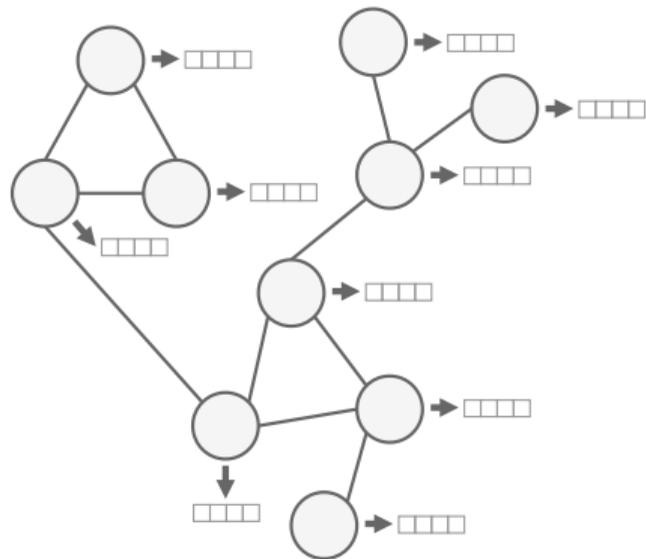
[25] J. Bruna et al., "Spectral networks and locally connected networks on graphs," *arXiv preprint arXiv:1312.6203*, 2013.

[26] M. M. Bronstein et al., "Geometric deep learning: going beyond euclidean data," *IEEE Signal Processing Magazine*, vol. 34, no. 4, 2017.

Graph-level v. node-level



Graph-level learning.
(e.g., molecules)



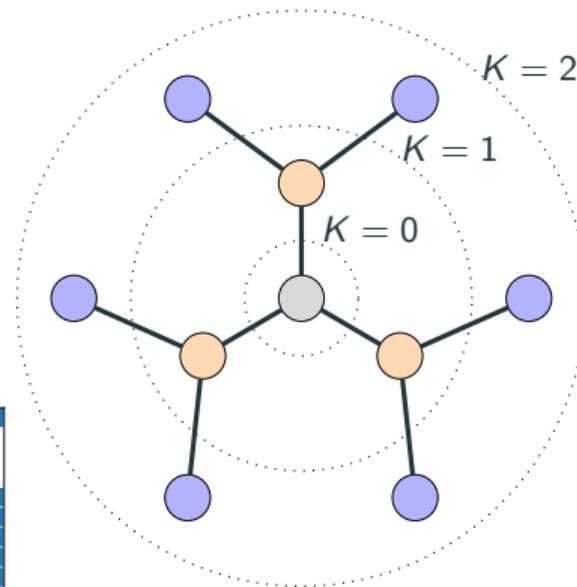
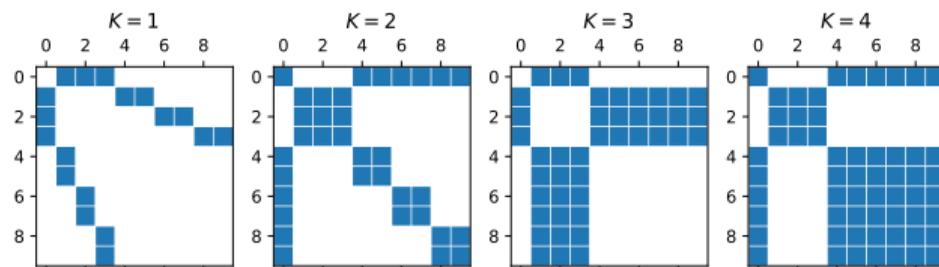
Node-level learning.
(e.g., social networks)

Powers of R

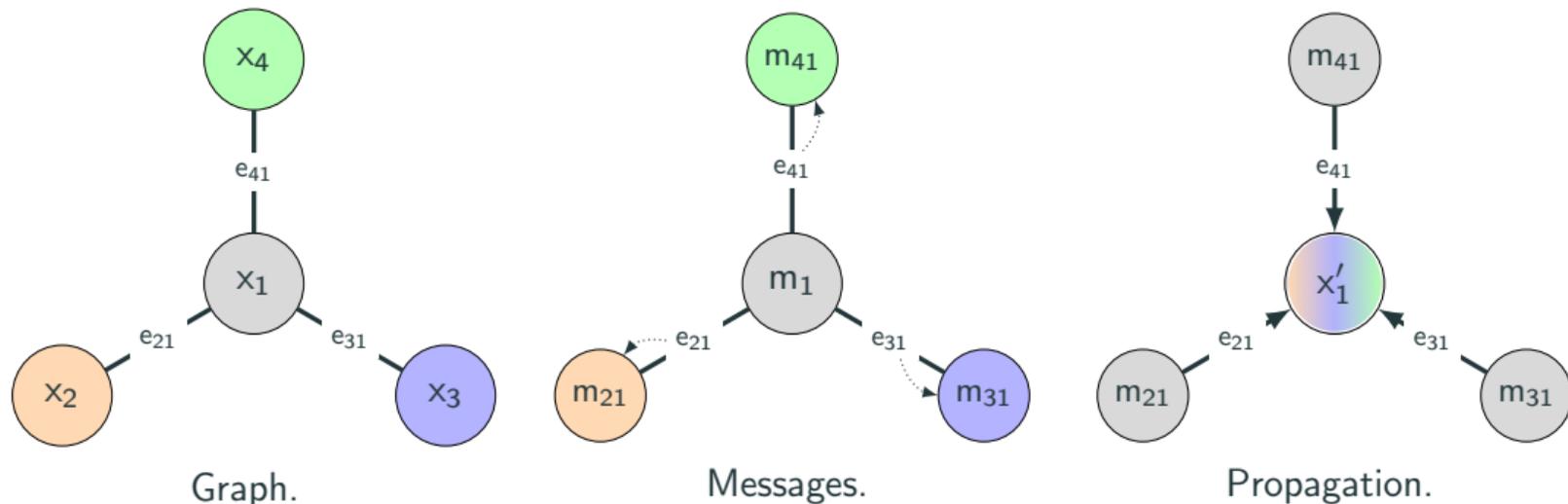
Let's consider the effect of applying R^2 to X :

$$(RRX)_i = \sum_{j \in \mathcal{N}(i)} r_{ij}(RX)_j = \sum_{j \in \mathcal{N}(i)} \sum_{k \in \mathcal{N}(j)} r_{ij} \cdot r_{jk} \cdot x_k$$

Key idea: by applying R^K we **read from the K -th order neighbourhood** of a node.



Message Passing Neural Networks [27]



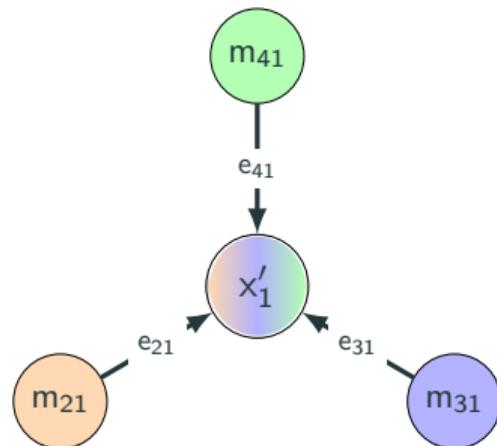
[27] J. Gilmer et al., "Neural message passing for quantum chemistry," *arXiv preprint arXiv:1704.01212*, 2017.

Message Passing Neural Networks [27]

A general scheme for message-passing networks:

$$x'_i = \gamma \left(x_i, \square_{j \in \mathcal{N}(i)} \phi(x_i, x_j, e_{ji}) \right),$$

- ϕ : **message function**, depends on x_i , x_j and possibly the edge attribute e_{ji} (we call messages m_{ji});
- $\square_{j \in \mathcal{N}(i)}$: **aggregation function** (sum, average, max, or something else...);
- γ : **update function**, final transformation to obtain new attributes after aggregating messages.



[27] J. Gilmer et al., "Neural message passing for quantum chemistry," *arXiv preprint arXiv:1704.01212*, 2017.

Chebyshev Polynomials [2]

A recursive definition using Chebyshev polynomials:

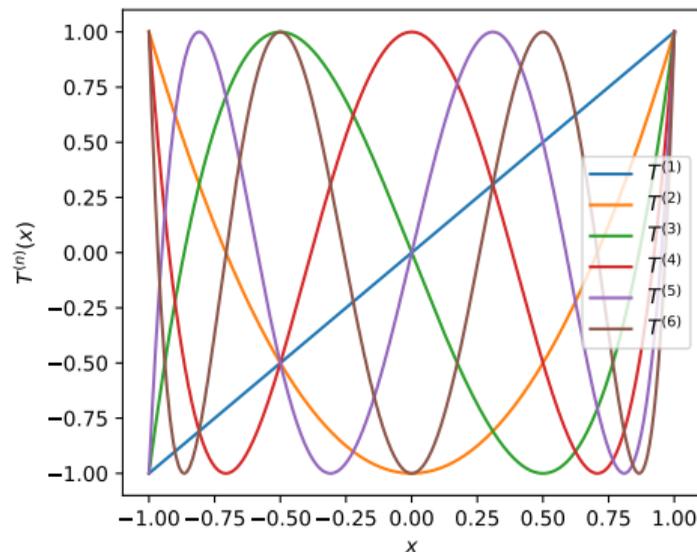
$$T^{(0)} = I$$

$$T^{(1)} = \tilde{L}$$

$$T^{(k)} = 2 \cdot \tilde{L} \cdot T^{(k-1)} - T^{(k-2)}$$

Where $\tilde{L} = \frac{2L_n}{\lambda_{\max}} - I$ and $L_n = I - D^{-1/2}AD^{-1/2}$

$$\text{Layer: } X' = \sigma \left(\sum_{k=0}^K T^{(k)} X \Theta^{(k)} \right)$$



[2] M. Defferrard et al., "Convolutional neural networks on graphs with fast localized spectral filtering," *Advances in Neural Information Processing Systems*, 2016.

Graph Convolutional Networks [3]

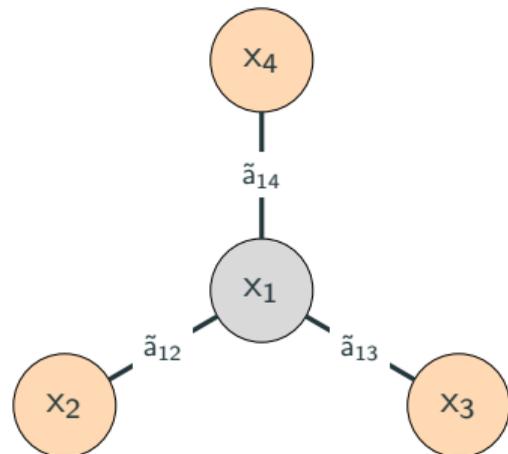
Polynomial of order $K \rightarrow K$ layers of order 1;

Three simplifications:

1. $\lambda_{\max} = 2 \rightarrow \tilde{L} = \frac{2L_n}{\lambda_{\max}} - I = -D^{-1/2}AD^{-1/2} = -A_n$
2. $K = 1 \rightarrow X' = X\Theta^{(0)} - A_nX\Theta^{(1)}$
3. $\Theta = \Theta^{(0)} = -\Theta^{(1)}$

$$\text{Layer: } X' = \sigma\left(\underbrace{(I + A_n)}_{\tilde{A}}X\Theta\right) = \sigma(\tilde{A}X\Theta)$$

$$\text{For stability: } \tilde{A} = D^{-1/2}(I + A)D^{-1/2}$$



[3] T. N. Kipf et al., "Semi-Supervised Classification with Graph Convolutional Networks," *International Conference on Learning Representations (ICLR)*, 2017.

Edge-Conditioned Convolution [28]

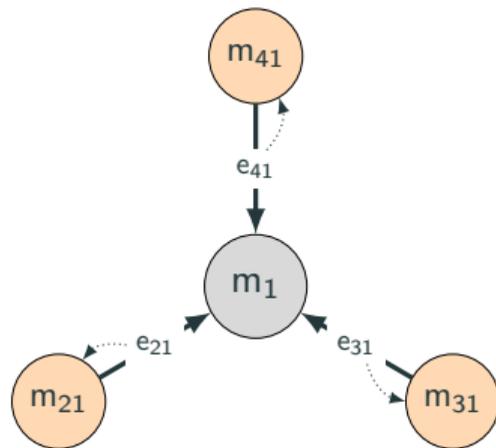
Key idea: incorporate edge attributes into the messages.

Consider a MLP $\phi : \mathbb{R}^S \rightarrow \mathbb{R}^{FF'}$ called a **filter generating network**:

$$\Theta^{(ji)} = \text{reshape}(\phi(e_{ji}))$$

Use the edge-dependent weights to compute messages:

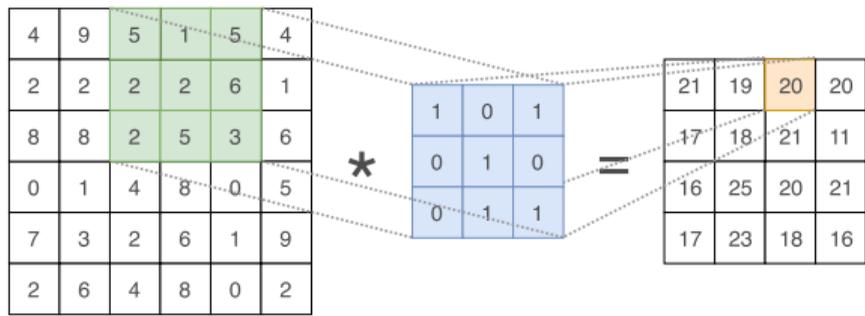
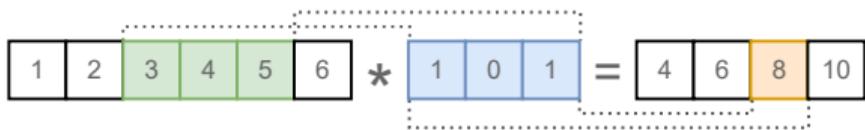
$$x'_i = \Theta^{(i)}x_i + \sum_{j \in \mathcal{N}(i)} \Theta^{(ji)}x_j + b$$



[28] M. Simonovsky et al., "Dynamic edgeconditioned filters in convolutional neural networks on graphs," 2017.

Graph Convolution

Discrete Convolution



Recall: CNNs compute a discrete convolution

$$(f \star g)[n] = \sum_{m=-M}^M f[n-m]g[m] \quad (1)$$

Convolution Theorem

Given two functions f and g , their convolution $f \star g$ can be expressed as:

$$f \star g = \mathcal{F}^{-1} \{ \mathcal{F} \{ f \} \cdot \mathcal{F} \{ g \} \} \quad (2)$$

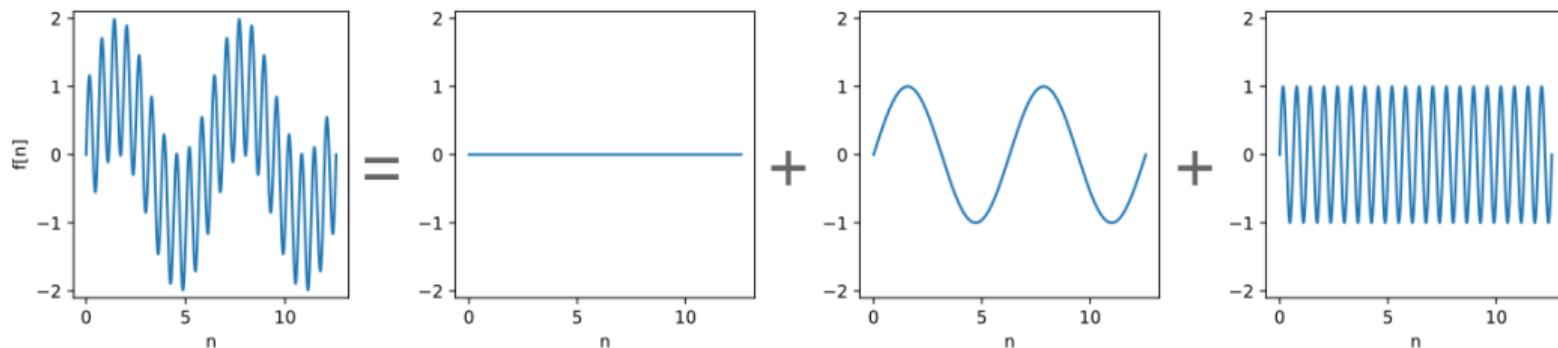
Where \mathcal{F} is the **Fourier transform** and \mathcal{F}^{-1} its inverse.

What is the Fourier transform?

Key intuition – we are representing a function in a different basis.

$$\mathcal{F}\{f\}[k] = \hat{f}[k] = \sum_{n=0}^{N-1} f[n] e^{-i \frac{2\pi}{N} kn}$$

$$\mathcal{F}^{-1}\{\hat{f}\}[n] = f[n] = \frac{1}{N} \sum_{k=0}^{N-1} \hat{f}[k] e^{i \frac{2\pi}{N} kn}$$

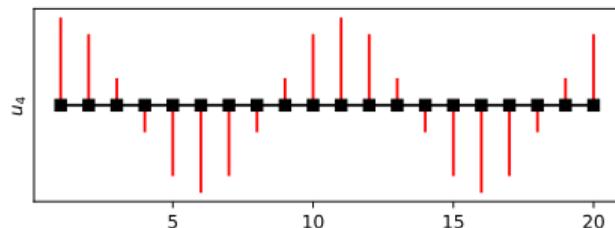
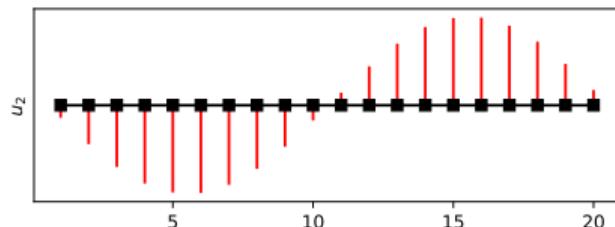
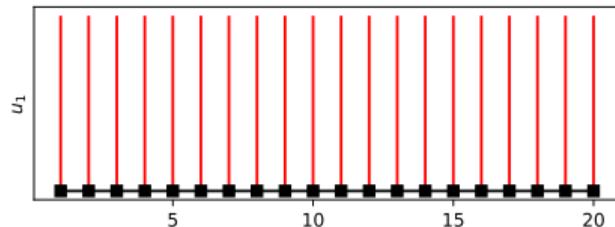


From FT to GFT

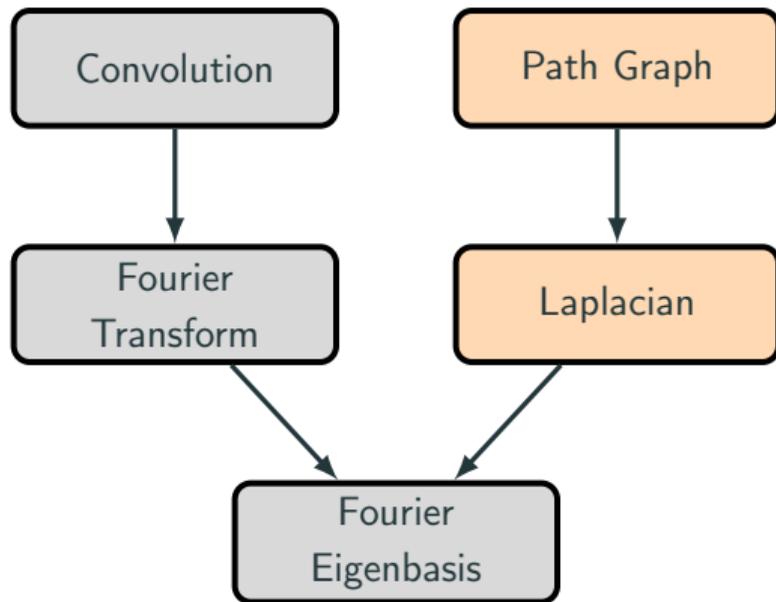
The eigenvectors of the Laplacian for a path graph can be obtained analytically:

$$u_k[n] = \begin{cases} 1, & \text{for } k = 0 \\ e^{i\pi(k+1)n/N}, & \text{for odd } k, k < N - 1 \\ e^{-i\pi kn/N}, & \text{for even } k, k > 0 \\ \cos(\pi n), & \text{for odd } k, k = N - 1 \end{cases}$$

Looks familiar?



From FT to GFT



- Drop the “grid” assumption
- Replace $e^{-i\frac{2\pi}{N}kn}$ with generic $u_k[n]$:

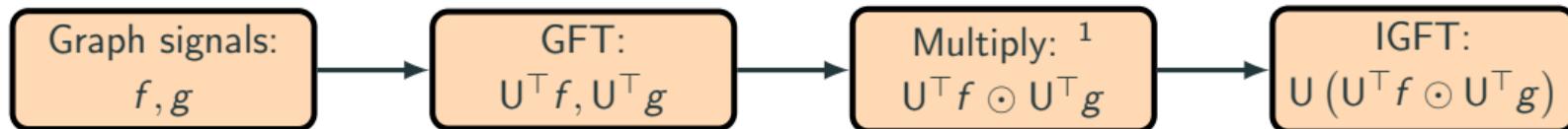
$$\mathcal{F}_G \{f\} [k] = \sum_{n=0}^{N-1} f[n]u_k[n]$$

- GFT: $\mathcal{F}_G \{f\} = \hat{f} = U^T f$;
- IGFT: $\mathcal{F}_G^{-1} \{\hat{f}\} = f = U\hat{f}$

Graph Convolution

Recall:

- Convolution theorem: $f \star g = \mathcal{F}^{-1} \{ \mathcal{F} \{ f \} \cdot \mathcal{F} \{ g \} \}$
- Spectral theorem: $L = U \Lambda U^T = \sum_{i=0}^{N-1} \lambda_i u_i u_i^T$

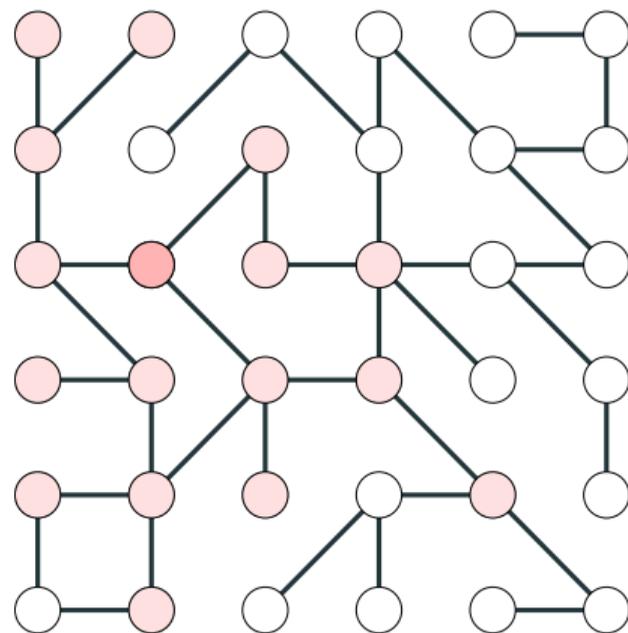
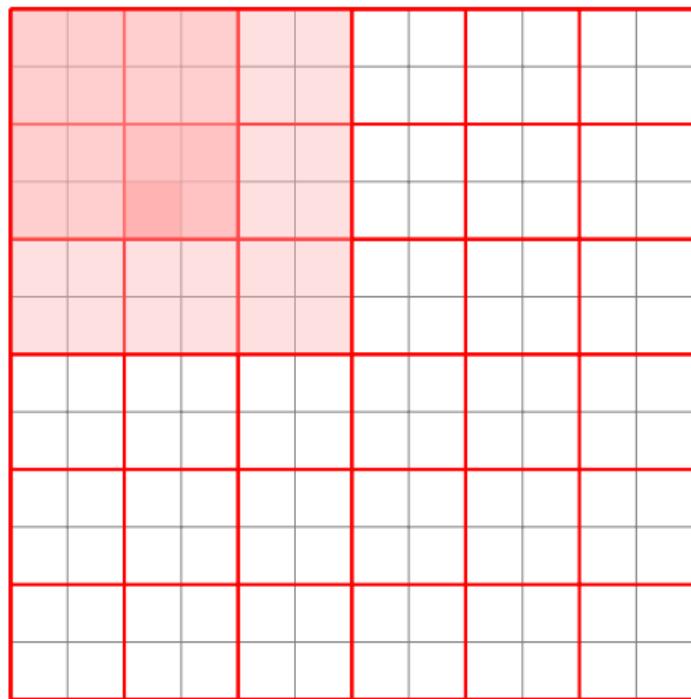


$$\text{Graph filter: } U (U^T f \odot U^T g) = U \cdot \underbrace{\text{diag}(U^T g)}_{g(\Lambda)} \cdot U^T f = \underbrace{U \cdot g(\Lambda)}_{g(L)} \cdot U^T f = g(L) f$$

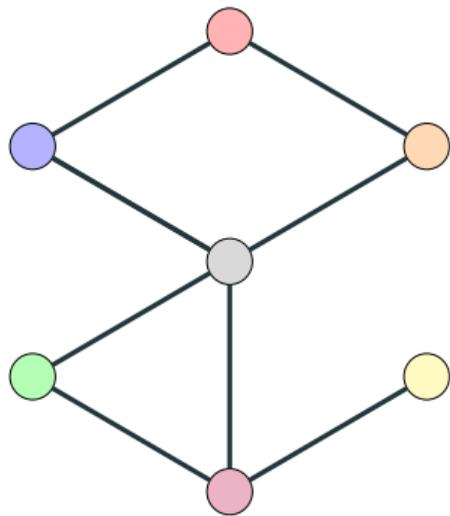
¹\odot indicates element-wise multiplication

Pooling operators

Pooling in CNNs



Selecting nodes



Example 1: partition.



Example 2: cover (possible overlaps).



Example 3: sparse.



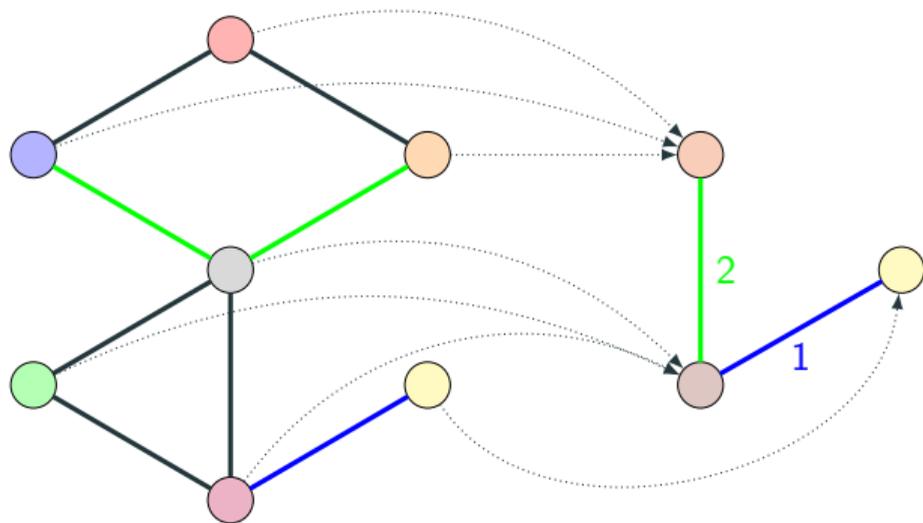
Select, Reduce, Connect [6]

Putting everything together:

$$\mathcal{S} = \underbrace{\{\mathcal{S}_k\}_{k=1:K}}_{\text{Selection}} = \text{Sel}(\mathcal{G});$$

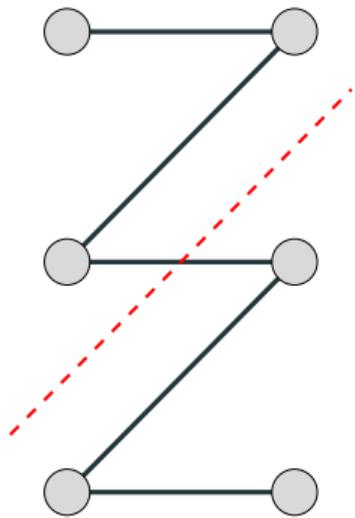
$$\mathcal{X}' = \underbrace{\{\text{Red}(\mathcal{G}, \mathcal{S}_k)\}_{k=1:K}}_{\text{Reduction}};$$

$$\mathcal{E}' = \underbrace{\{\text{Con}(\mathcal{G}, \mathcal{S}_k, \mathcal{S}_l)\}_{k,l=1:K}}_{\text{Connection}};$$

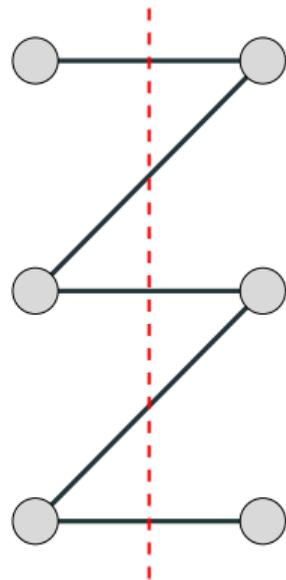


[6] D. Grattarola et al., "Understanding Pooling in Graph Neural Networks," *Under review at IEEE TNMLS*, 2021.

Mincut vs. Maxcut

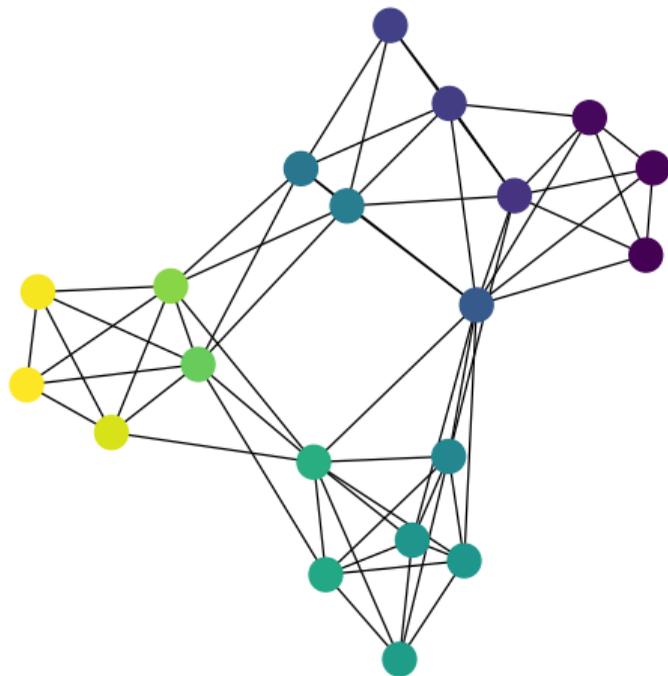


MinCut

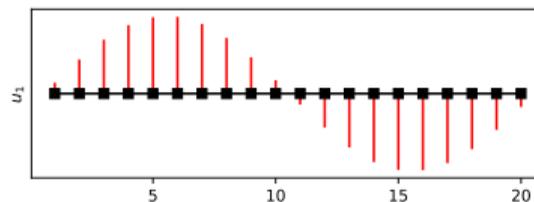


MaxCut

Spectral clustering [30]



The low-frequency eigenvectors naturally cluster the nodes.

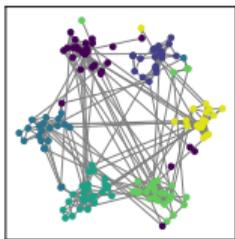


Idea: run k-means clustering (or similar) using the first few eigenvectors.

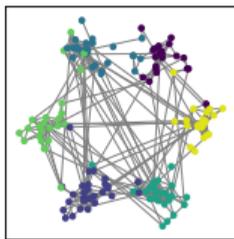
[29] J. Shi et al., "Normalized cuts and image segmentation," *Departmental Papers (CIS)*, 2000.

[30] U. Von Luxburg, "A tutorial on spectral clustering," *Statistics and computing*, vol. 17, no. 4, 2007.

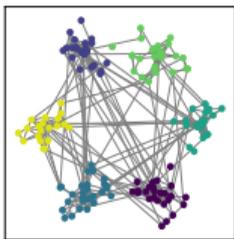
MinCut pooling



(a) SC



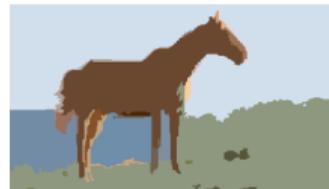
(b) DiffPool



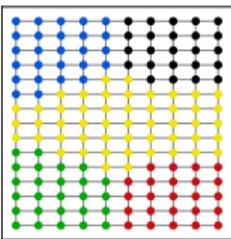
(c) MinCut



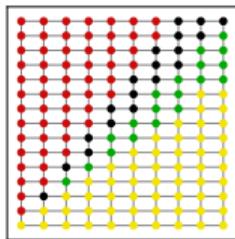
(a) Original



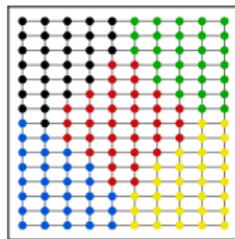
(b) SC



(d) SC



(e) DiffPool



(f) MinCut

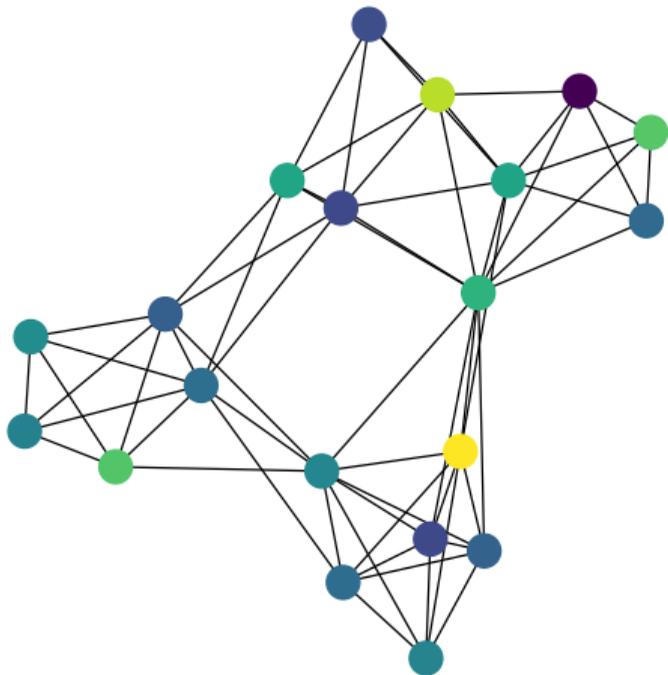


(c) DiffPool

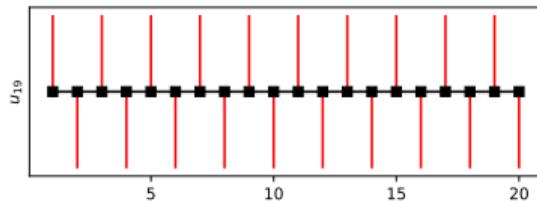


(d) MinCut

Node decimation [13]



Alternative: use the highest-frequency eigenvector to do something similar to a regular subsampling.



[31] L. Palagi et al., "Computational approaches to max-cut," 2012.

[13] F. M. Bianchi et al., "Hierarchical representation learning in graph neural networks with node decimation pooling," *IEEE Transactions on Neural Networks and Learning Systems*, 2020.

Learning to pool

Key idea: learn to output S^T by giving node features X as input to a neural network.

- **DiffPool** [7]: GNN for S^T , regularize with “link prediction” loss;
- **MinCutPool** [8]: MLP for S^T , regularize with “minimum cut” loss (same objective as spectral clustering);
- **Deep Graph Mapper** [32]: combine Mapper [33] and GCN [3] to compute clusters.

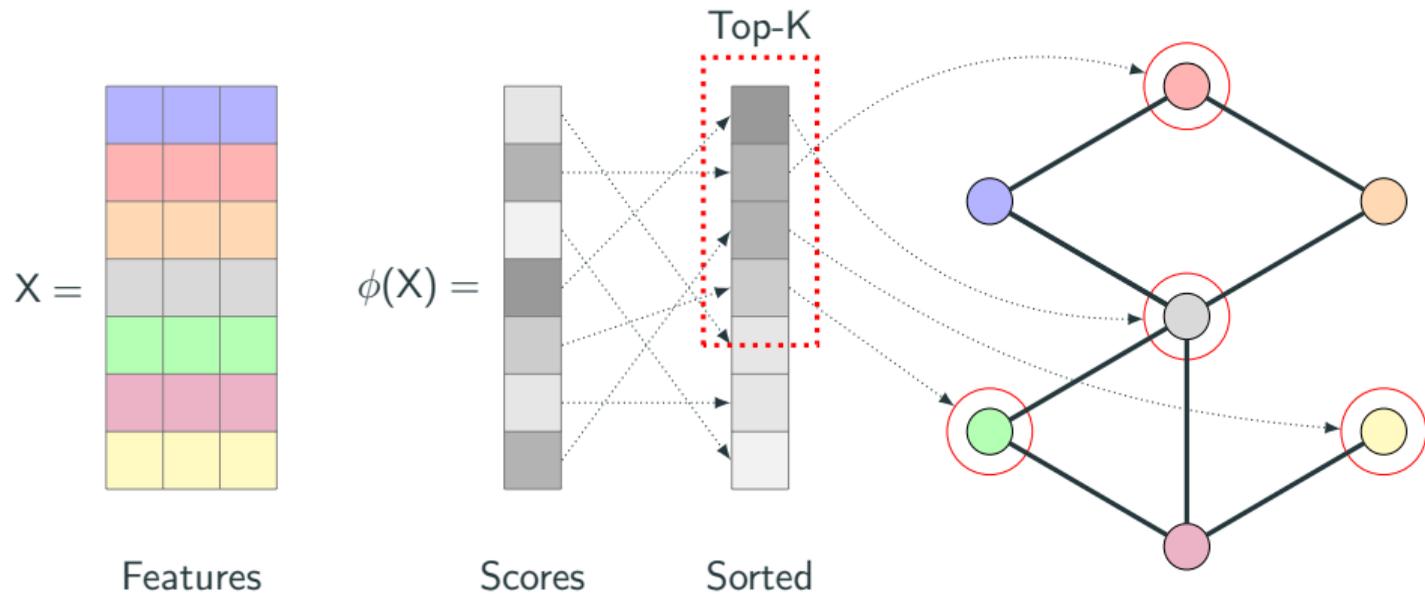
$$\phi(X) = S^T = \begin{array}{|c|c|c|} \hline \text{blue} & \text{light blue} & \text{light blue} \\ \hline \text{red} & \text{light red} & \text{light red} \\ \hline \text{orange} & \text{light orange} & \text{light orange} \\ \hline \text{grey} & \text{light grey} & \text{light grey} \\ \hline \text{green} & \text{light green} & \text{light green} \\ \hline \text{pink} & \text{light pink} & \text{light pink} \\ \hline \text{yellow} & \text{light yellow} & \text{yellow} \\ \hline \end{array} \in \mathbb{R}^{N \times K}$$

[7] R. Ying et al., “Hierarchical Graph Representation Learning with Differentiable Pooling,” *arXiv preprint arXiv:1806.08804*, 2018.

[8] F. M. Bianchi et al., “Spectral Clustering with Graph Neural Networks for Graph Pooling,” *International Conference on Machine Learning*, 2020.

[32] C. Bodnar et al., “Deep Graph Mapper: Seeing Graphs through the Neural Lens,” *arXiv preprint arXiv:2002.03864*, 2020.

Top-K methods



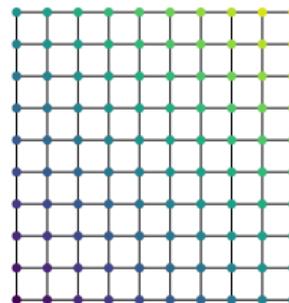
Top-K methods

Reduce: $X' = X_i$ - Connect: $A' = A_{i,i}$

Problems:

- Top-k selection is **non-differentiable** (no way of training ϕ). Solved by **gating** (multiplying) the node attributes with the scores.
- Graph is likely to be **disconnected** or simply **cut off** (like in the image on the right).
Not really solvable...

Original



Top-K



Main properties of pooling operators

- **Dense vs. Sparse**: how many nodes are selected for the supernodes;
- **Fixed vs. Adaptive**: how many supernodes does the selection compute;
- **Trainable vs. Non-trainable**: learn to pool from data or not;

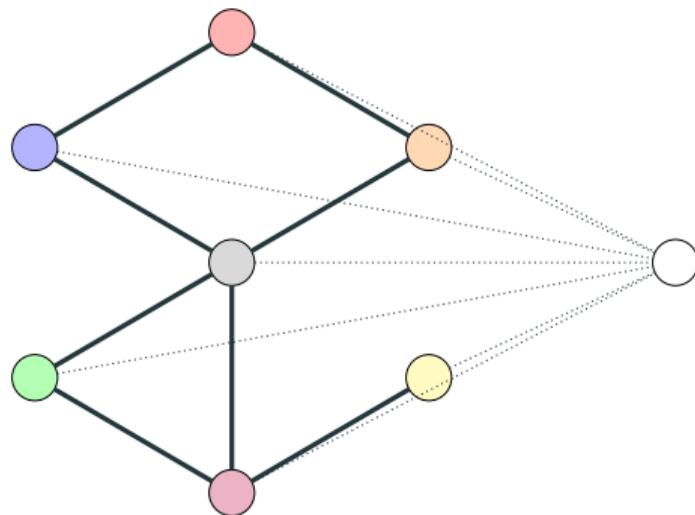
Global Pooling

In CNNs, after convolution, we usually **flatten** out the images to give a vector as input to a MLP:

1	2	3
4	5	6
7	8	9

1	2	3	4	5	6	7	8	9
---	---	---	---	---	---	---	---	---

The graph equivalent must be **invariant to permutations** of the nodes:



More results

Molecule generation [17]

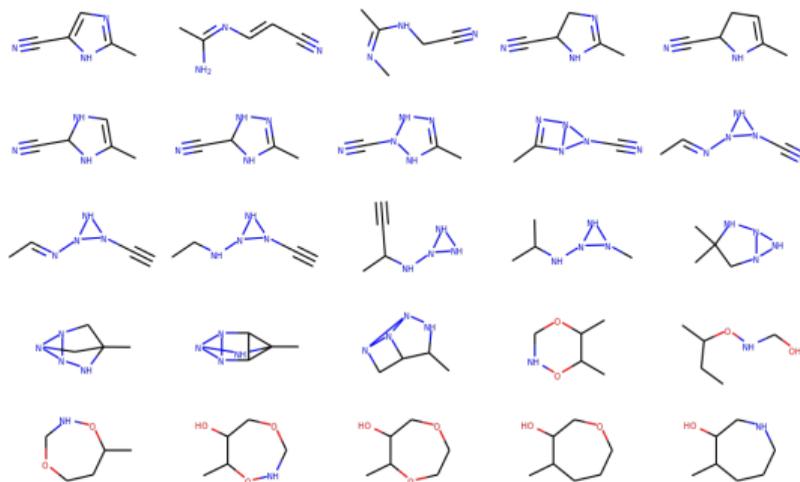
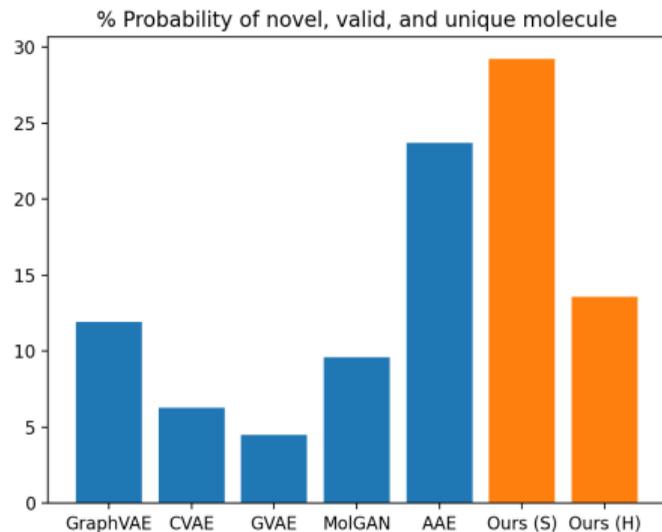


Figure 8: Valid, novel, and unique molecules

[17] D. Grattarola et al., "Adversarial autoencoders with constant-curvature latent manifolds," *Applied Soft Computing*, 2019.