

Framelets and Graph Neural Networks

Yu Guang Wang

SJTU INS & Math

yuguang.wang@sjtu.edu.cn

SJTU Math Salon 2021



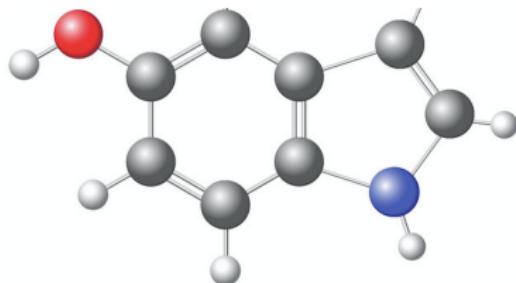
上海交通大学
SHANGHAI JIAO TONG UNIVERSITY

What are Graph Neural Networks?

Graph neural networks (GNNs) are **deep nets** which take **graph-structure data** as input.

- The structure of a typical GNN is similar as CNN, which has **multiple neural layers**, each of which contains a number of neuronal nodes.
- The input of GNN is graph-structured data which usually consists of **adjacency matrix** (representing edges) plus **features on vertices**.

Inputs of GNN and CNN



Graph data for GNN

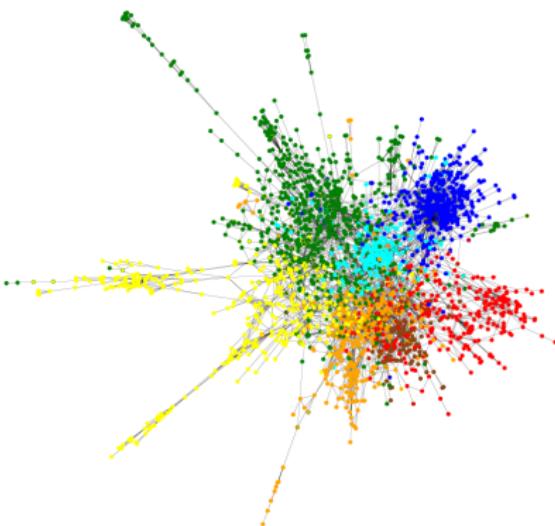
nodes, edges
size varies



Image for CNN

pixels
regular grid

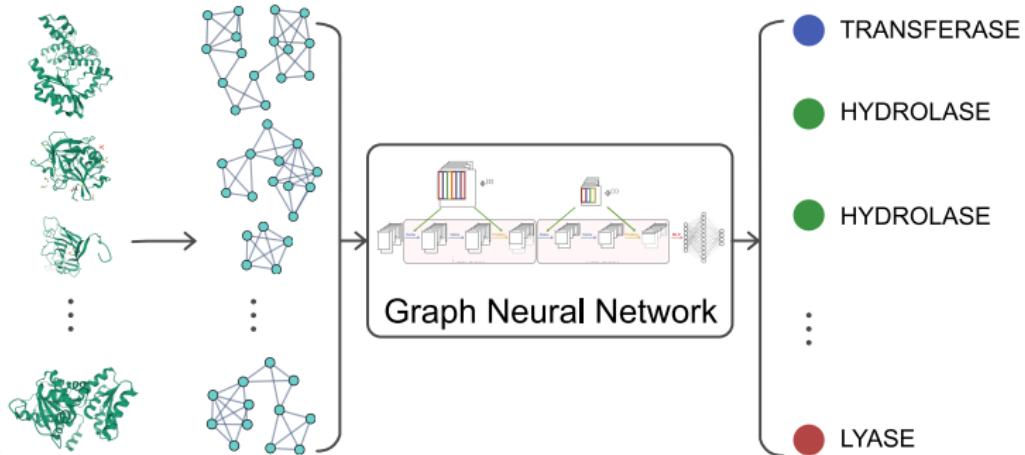
What problems can Graph Neural Networks solve?



Node property prediction (node classification, drug repositioning (knowledge graph), recommender systems e.g. e-commerce)

- The input is one graph.
- The GNN uses the known labeled nodes and the graph edges to infer the missing labels.

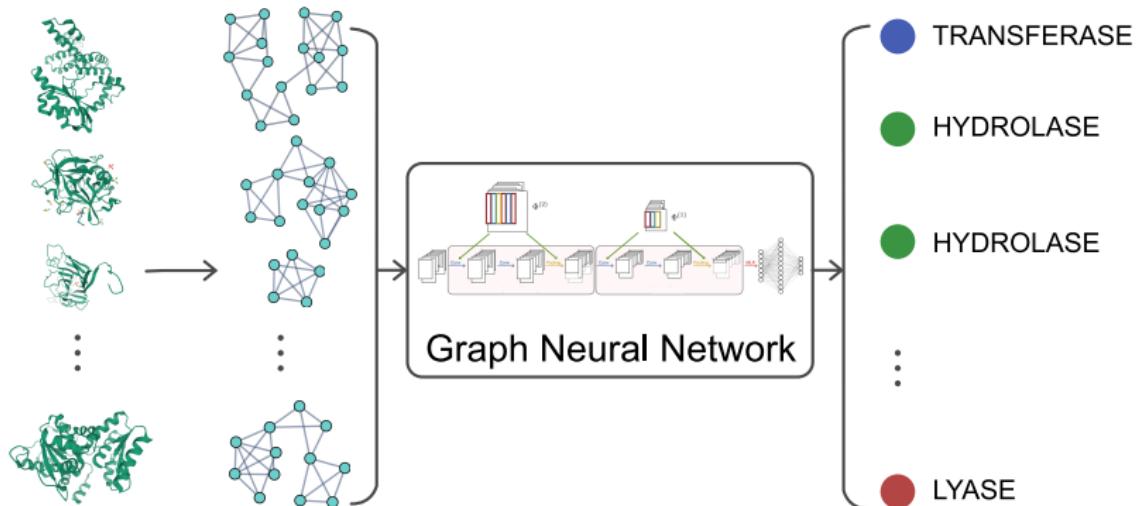
What problems can Graph Neural Networks solve? (Continued)



Graph property prediction (graph classification/regression, graph generation, 3D object recognition e.g. LiDAR in self-driving)

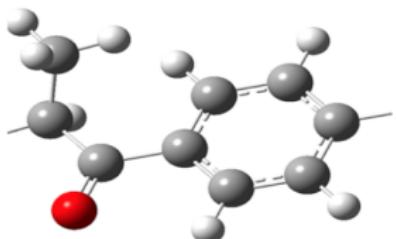
- The input is a set of graphs. Each graph sample has a Y-label.
- The graph samples with known Y-labels and features on the graphs are used in training for GNN. The trained GNN model can be used to infer the missing labels.

Graph Classification for Enzymes



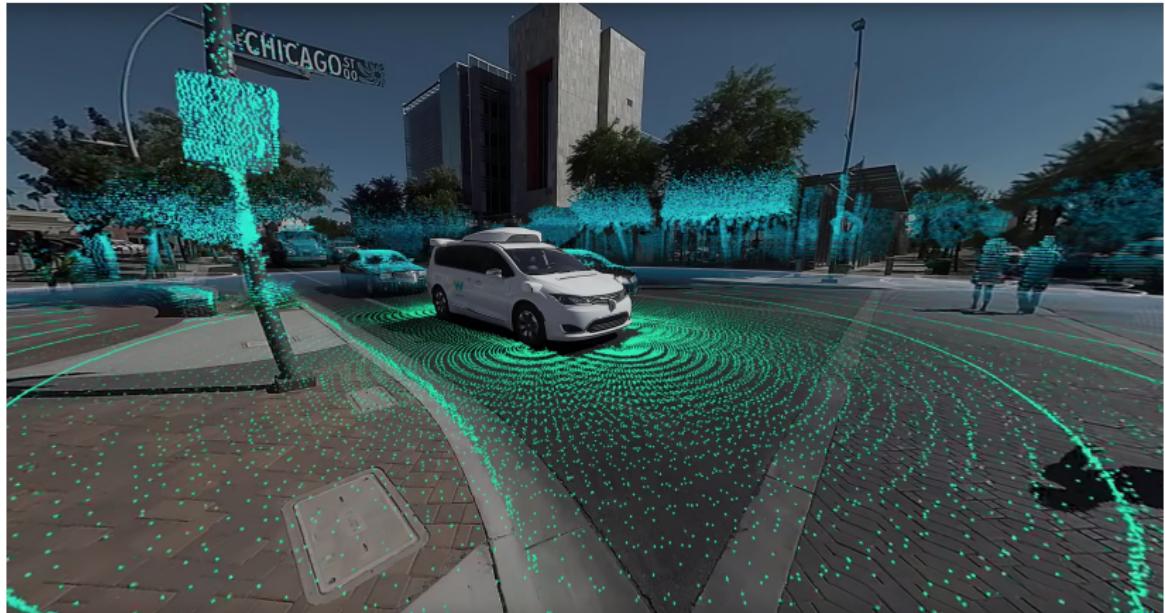
GNNs for labeling enzymes. A set of graphs is extracted as the input data set, where each graph represents a specific protein tertiary structure. The task is to assign each enzyme instance to one of the given EC top-level classes correctly. The GNN's role is to find universally applicable rules to label the graphs by learning the topological and feature information of the input. The structure of enzymes is retrieved from the Protein Data Bank.

Quantum Chemistry Graph Regression



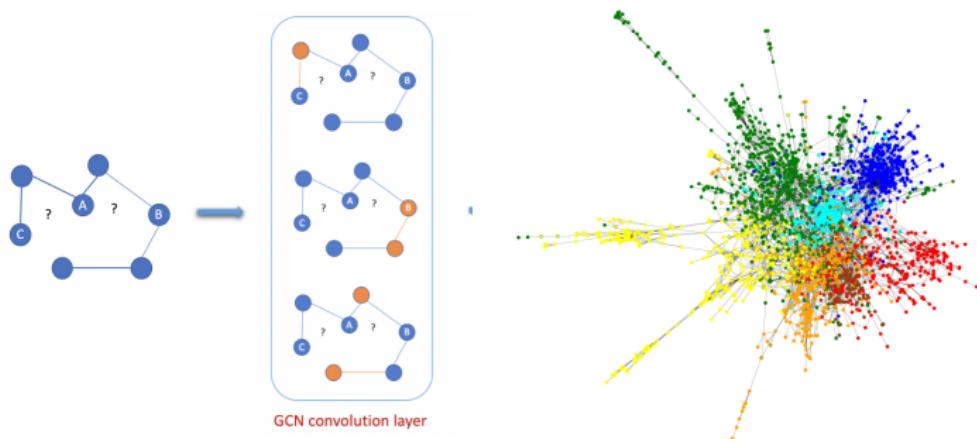
- QM7 is a collection of 7,165 molecules, train/test = 4/1.
- Each molecule contains ≤ 23 atoms (including C, O, N, S), atoms are connected by bonds, molecular structure varies (e.g. double/triple bonds, cycles, carboxy, cyanide ...).
- Molecule is a graph, atoms are nodes, bonds are edges and Coulomb energy as weights, then Coulomb energy matrix is adjacency matrix.
- Task: to predict atomization energy of molecule given the molecular structure.

3D Object Recognition for Self-driving LiDAR



Source: Waymo

What makes Graph Neural Networks function?



Spectral graph convolution $g \star f = U((U^T g) \odot (U^T f))$

- Motivated by traditional Fourier convolution
- Have similar role as convolution in CNN
- Can preserve the structural information of graph data input.
- More general form of neural message passing *Gilmer et al. (ICML 2017)*.

Framelet convolution

For network filter θ and input graph feature $X \in \mathbb{R}^{N \times d}$ of the graph \mathcal{G} with N nodes, we define

$$\theta \star X = \text{ReLU}(\mathcal{V}(\text{diag}(\theta)(WX'))), \quad X' = XW.$$

- WX' is the framelet coefficient matrix for the transformed X'
- \mathcal{W} is a sequence of $nJ + 1$ transform matrices (each of size $N \times N$) for low-pass and high-passes
- The size of the vector θ is $(nJ + 1)N$ which matches the total number of the framelet coefficients for each feature.
- The network filter θ lies in the frequency domain, each component of which is multiplied to the corresponding row of WX' .
- The matrix W is a trainable weight matrix with dimension $d \times d'$.

Graph framelets

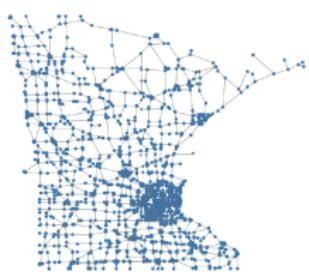
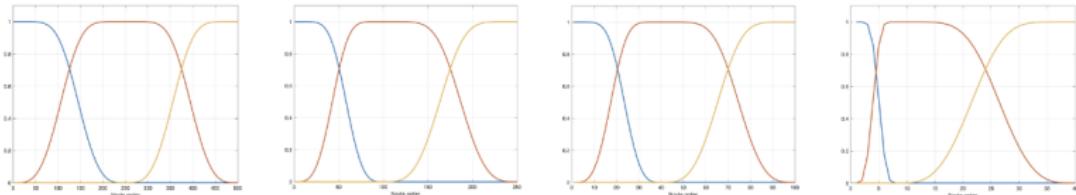
Suppose $\{(\lambda_\ell, \mathbf{u}_\ell)\}_{j=1}^N$ are the eigenvalue and eigenvector pairs for \mathcal{L} of graph \mathcal{G} with N nodes. The undecimated framelets at scale level $j = 1, \dots, J$ for graph \mathcal{G} , for $n = 1, \dots, r$, by

$$\varphi_{j,p}(v) = \sum_{\ell=1}^N \widehat{\alpha}\left(\frac{\lambda_\ell}{2^j}\right) \overline{\mathbf{u}_\ell(p)} \mathbf{u}_\ell(v)$$
$$\psi_{j,p}^n(v) = \sum_{\ell=1}^N \widehat{\beta^{(n)}}\left(\frac{\lambda_\ell}{2^j}\right) \overline{\mathbf{u}_\ell(p)} \mathbf{u}_\ell(v)$$

- $\varphi_{j,p}$ or $\psi_{j,p}^r$ is the low-pass or high-pass framelet translated at node p .
- The low-pass and high-pass framelet coefficients for a signal f on graph \mathcal{G} are $v_{j,p}$ and $w_{j,p}^r$, which are the projections $\langle \varphi_{j,p}, f \rangle$ and $\langle \psi_{j,p}^r, f \rangle$ of the graph signal onto framelets at scale j and node p .
- The dilation factor is 2^j with the dilation (base) 2.

Framelet transforms via filter bank (Minnesota traffic network)

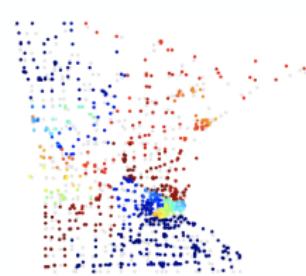
Filter bank with 2 high passes



(a) Original Network



(c) \hat{w}_0^1



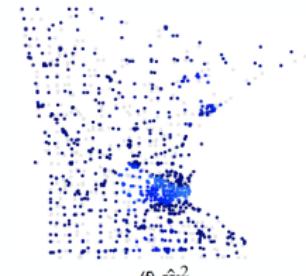
(e) \hat{w}_1^1



(b) \hat{v}_0

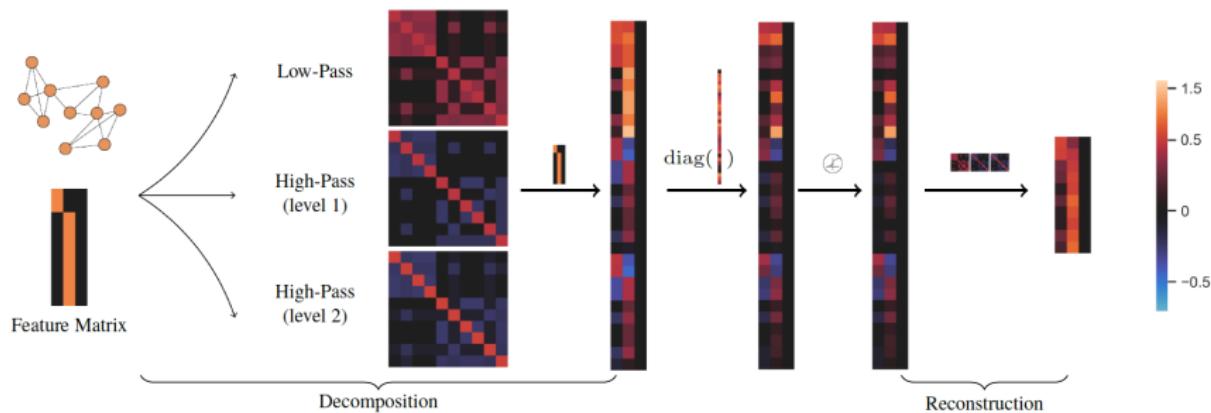


(d) \hat{w}_0^2



(f) \hat{w}_1^2

Shrinkage activation



$$\boldsymbol{\theta} \star X = V \left(\text{Shrinkage} \left(\text{diag}(\boldsymbol{\theta})(WX') \right) \right), \quad X' = XW.$$

- Different from ReLU, the shrinkage activation works in the framelet domain when the shrinkage thresholds for the high-pass coefficients in the framelet domain.

Shrinkage for signal compression

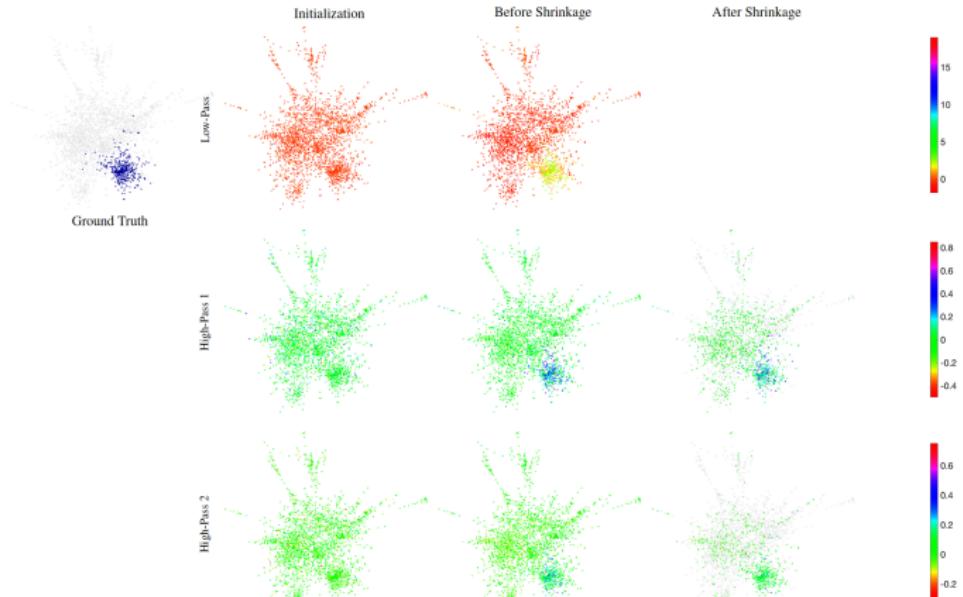
The high-pass coefficients in the frequency domain can be cut off by shrinkage thresholding. For example, the *soft-thresholding* *Donoho (1994), Donoho (1995), Tibshirani (1996)*. defines

$$\text{Shrinakge}(x) = \text{sgn}(x)(|x| - \lambda)_+ \quad \forall x \in \mathbb{R},$$

where λ is the threshold value.

- Any x with its absolute value less than λ shall return to zero.
- Applying the above soft-thresholding to the shrinkage activation in framelet convolution only influences small high-pass framelet coefficients.
- We also consider the scale-dependent selection threshold with demarcation point *Donoho (1995)*. : $\lambda = \sigma \sqrt{2 \log(N)} / \sqrt{N}$ for N coefficients.
- The hyperparameter σ is an analogue to the noise level of the wavelet denoising model. Here we let σ be associated with the magnitude order of the coefficients so it reflects the scale of the framelet representation.

Framelet convolution on Cora



- Coefficients at initialization, after 2 convolutional layers before shrinkage activation, and after shrinkage activation. (Comp Ratio 47.7%)
- Horizontal comparison indicates that high-pass coefficients compressed a critical part of coefficients after shrinkage (in green).
- Vertical comparison shows high-pass coefficients usually have more distinctive values and concentrated on the detailed information compared to low-pass.

Node classification benchmark

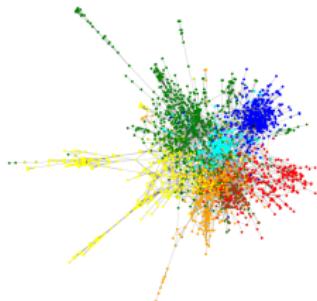
The node classification tasks for GNN are conducted on Cora, Citeseer and Pubmed, which are three benchmark citation networks. Moreover, we employ ogbn-arxiv from open graph benchmark OGB to illustrate the power of our framelet convolution on large-scale graph-structured data.

Table 5. Summary of the datasets for node classification tasks.

	Cora	Citeseer	Pubmed	ogbn-arxiv
# Nodes	2,708	3,327	19,717	169,343
# Edges	5,429	4,732	44,338	1,166,243
# Features	1,433	3,703	500	128
# Classes	7	6	3	40
# Training Nodes	140	120	60	90,941
# Validation Nodes	500	500	500	29,799
# Test Nodes	1,000	1,000	1,000	48,603
Label Rate	0.052	0.036	0.003	0.537

- The UFGConv-R and UFGConv-S are compared against MLP, DeepWalk, Chebyshev, GCN, Spectral CNN, GWNN, MPNN, GraphSAGE, LanczosNet, DCNN, GAT.

Framelet convolution for Node classification

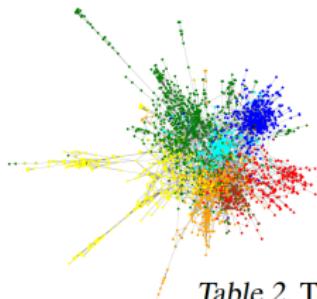


We report the accuracy score in percentage with the top-3 highlighted in the following Table 2. For UFGCONV-S, we also report the compression ratio for shrinkage (in green). The UFGCONV-R method achieves the highest prediction accuracy among all baseline models. The learned UFGCONV-S with threshold level $\sigma = 1$ trims up to 50% information but still obtains the top-3 rank on two tasks.

Method	Cora	Citeseer	Pubmed
MLP	55.1	46.5	71.4
DEEPWALK	67.2	43.2	65.3
SPECTRAL	73.3	58.9	73.9
CHEBYSHEV	81.2	69.8	74.4
GCN	81.5	70.3	79.0
GWNN	82.8	71.7±0.6	79.1
GAT	83.0±0.7	72.5±0.7	79.0±0.3
MPNN	78.0±1.1	64.0±1.9	75.6±1.0
GRAPHSAGE	74.5±0.8	67.2±1.0	76.8±0.6
LANCZOSNET	79.5±1.8	66.2±1.9	78.3±0.3
DCNN	79.7±0.8	69.4±1.3	76.8±0.8
UFGCONV-S (Compression)	83.0±0.5 (47.7)	71.0 (39.0)	79.4±0.4 (27.7)
UFGCONV-R	83.6±0.6	72.7±0.6	79.9±0.1

† The top three are highlighted by **First**, **Second**, **Third**.

Framelet convolution for ogbn-arXiv (169K nodes, 1.1m edges)



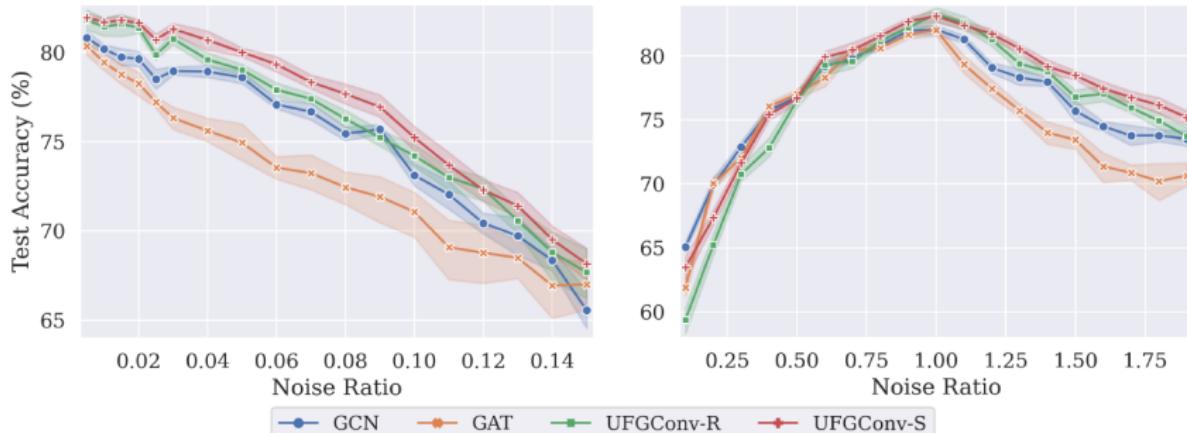
Framelet convolution achieves similar outstanding performance as reported in Table 2 for the **ogbn-arxiv** dataset, where the **UFGConv-R** ranks first with a moderate number of parameters, and the **UFGConv-S** with threshold $\sigma = 1$ achieves a comparably excellent accuracy using **64.2%** information.

Table 2. Test accuracy (in percentage) for **ogbn-arxiv** with standard deviation after \pm . The compression ratio for UFGConv-S with shrinkage threshold level $\sigma = 1$ is **64.2%**.

Method	Test Acc.	Val. Acc.	#Params
MLP	55.50 \pm 0.23	57.65 \pm 0.12	110,120
NODE2VEC	70.07 \pm 0.13	71.29 \pm 0.13	21,818,792
GRAPHZOOM	71.18 \pm 0.18	72.20 \pm 0.07	8,963,624
P&L + C&S	71.26 \pm 0.01	73.00 \pm 0.01	5,160
GRAPHSAGE	71.49 \pm 0.27	72.77 \pm 0.17	218,664
GCN	71.74 \pm 0.29	73.00 \pm 0.17	142,888
DEEPERGCN	71.92 \pm 0.17	72.62 \pm 0.14	491,176
SIGN	71.95 \pm 0.11	73.23 \pm 0.06	3,566,128
GAAN	71.97 \pm 0.18	—	1,471,506
UFGConv-S	70.04 \pm 0.22	71.04 \pm 0.11	1,633,183
UFGConv-R	71.97 \pm 0.12	73.21 \pm 0.05	1,633,183

† The top three are highlighted by **First**, **Second**, **Third**.

Robustness of framelet convolution under noise

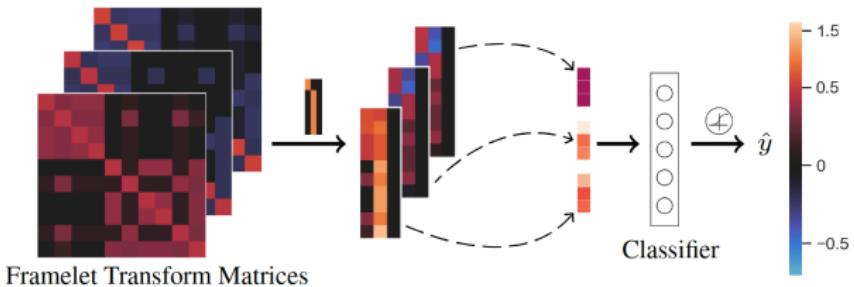


The left and right pictures show node and structure perturbation analysis on **Cora** respectively. We randomly change node feature or edge weight from 0/1 to its opposite 1/0. Figure reports the test accuracy of UFGConv, GCN and GAT. The x -axis is the distortion ratio (or SNR). As the SNR increases, UFGConv behaves well ahead of GCN and GAT while with higher test accuracy and smaller variance. The slightly lower performance of UFGConv only occurs when misinformation dominates the graph structure. It illustrates the effectiveness of shrinkage framelet convolution node and structure denoising.

How GNNs handle input graphs with varying number of nodes and connectivity structures?

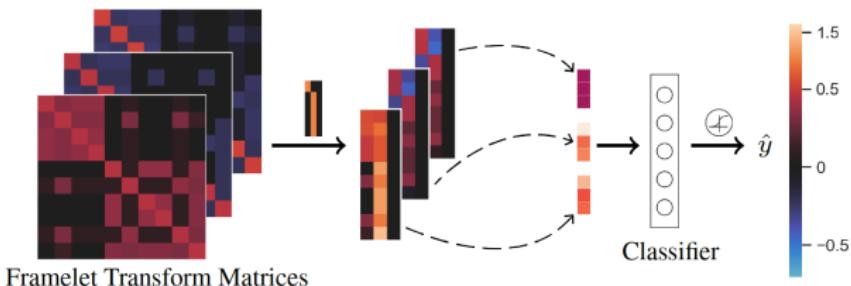
- One way is to use graph pooling. It is a computational strategy to reduce the number of graph nodes while preserve as much as geometric information of the original input graph data.
- By pooling, one has a unified graph-level representation for graph-structured data when the size and topology of an individual graph are changing.

Framelet pooling



- We use **2-level framelet system**. Given a graph with feature matrix $X \in \mathbb{R}^{N \times d}$, we can obtain a set of framelet coefficients $\mathcal{W}_{r,j}f$ including one low pass $\mathcal{W}_{0,2}f$ at level 2 and two high passes $\mathcal{W}_{1,1}f$ and $\mathcal{W}_{1,2}f$ at levels 1 and 2, respectively.
- Each scale-wise framelet coefficient is an $N \times d$ real-valued matrix, and its i th feature column $(\mathcal{W}_{r,j}X)_i$ for $i = 1, \dots, d$ would be aggregated by sum, or sum of squares of the elements.
- The two aggregation methods correspond to two *framelet pooling* strategies as shown in the picture. The calculation compresses the $N \times d$ coefficients to a d -dimensional vector, and the **pooled output** from the three framelet coefficients results in $3D$ vectors.

Energy conservation in framelet pooling



- The framelet pooling benefits the network training by employing the information from multi-scales, when all scales in the framelet representation of the graph signal are taken into account. Depending on how we aggregate the framelet coefficients, we distinguish the strategies by **UFGPOOL-SUM** or **UFGPOOL-SPECTRUM**.
- The latter aggregates the nodes by the wavelet (power) spectrum (i.e., the sum of absolute squares of framelet coefficients over nodes, $\sum_{p \in V} |v_{j,p}|^2$ and $\sum_{p \in V} |w_{j,p}^r|^2$). The total information of the graph signal X^{in} is then well-conserved after the pooling. The sum of wavelet power spectrum is equal to the total energy of the signal: $\|X^{\text{pooled}}\| = \|X^{\text{in}}\|$.

Framelet pooling for graph classification

We select six benchmarks to test the proposed pooling strategies, including four graph classification tasks with moderate sample sizes, one regression task, and one large-scale classification task. First five tasks use **TUDataset benchmarks**, including **D&D**, **PROTEINS** to categorize proteins into enzyme and non-enzyme structures; **NCI1** to identify chemical compounds that block lung cancer cells; **Mutagenicity** to recognize mutagenic molecular compounds for potentially marketable drug; and **QM7** to predict atomization energy value of molecules. The dataset **ogbg-molhiv** is for large-scale molecule classification.

Table 6. Summary of the datasets for the graph property prediction tasks.

Datasets	PROTEINS	Mutagenicity	D&D	NCI1	ogbg-molhiv	QM7
# Graphs	1,113	4,337	1,178	4,110	41,127	7,165
Min # Nodes	4	4	30	3	2	4
Max # Nodes	620	417	5,748	111	222	23
Avg # Nodes	39	30	284	30	26	15
Avg # Edges	73	31	716	32	28	123
# Features	3	14	89	37	9	0
# Classes	2	2	2	2	2	1 (R)

- We compare our framelet pooling **UFGPool-SUM** and **UFGPool-SPECTRUM** with six baseline methods that are capable for global pooling to verify the effectiveness of the learned graph representation. The baselines include **TopKPool**, **AttentionPool**, **SAGPool**, and the classic **Sum**, **Mean** and **Max** pooling.

Framelet pooling for Graph classification

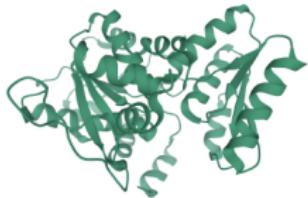


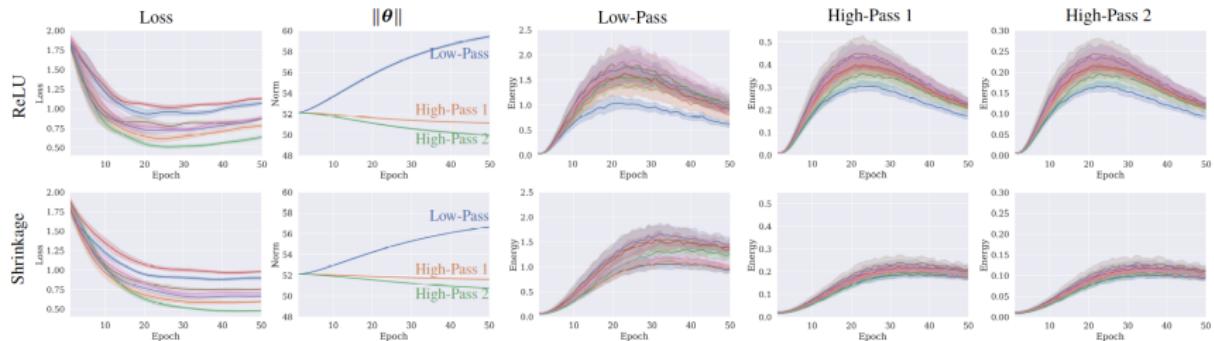
Table 3 shows our **UFGPOOL** methods outperform other methods on all datasets. Specifically, **UFGPOOL-SUM** achieves the top accuracy in four out of six datasets, and the second best accuracy in the other two, where the top performance is achieved by **UFGPOOL-SPECTRUM**. We also observe that **UFGPOOL-SPECTRUM** performs better on small molecules prediction: **Mutagenicity**, **QM7** and **ogbg-molhiv**. This precedence might come from encoding the multi-scale signal energy to the network where the **framelet spectra** capture the practically significant features of molecular data.

Table 3. Performance comparison for graph property prediction. **QM7** is a regression task in MSE; **ogbg-molhiv** is a classification task in ROC-AUC in percentage; others are for classification in test accuracy in percentage. The value after \pm is standard deviation.

Datasets	PROTEINS	Mutagenicity	D&D	NCI1	ogbg-molhiv	QM7
TOPKPOOL	73.48 \pm 3.57	79.84 \pm 2.46	74.87 \pm 4.12	75.11 \pm 3.45	78.14 \pm 0.62	175.41 \pm 3.16
ATTENTION	73.93 \pm 5.37	80.25 \pm 2.22	77.48 \pm 2.65	74.04 \pm 1.27	74.44 \pm 2.12	177.99 \pm 2.22
SAGPOOL	75.89 \pm 2.91	79.86 \pm 2.36	74.96 \pm 3.60	76.30 \pm 1.53	75.26 \pm 2.29	41.93 \pm 1.14
SUM	74.91 \pm 4.08	80.69 \pm 3.26	78.91 \pm 3.37	76.96 \pm 1.70	77.41 \pm 1.16	42.09 \pm 0.91
MAX	73.57 \pm 3.94	78.83 \pm 1.70	75.80 \pm 4.11	75.96 \pm 1.82	78.16 \pm 1.33	177.48 \pm 4.70
MEAN	73.13 \pm 3.18	80.37 \pm 2.44	76.89 \pm 2.23	73.70 \pm 2.55	78.21 \pm 0.90	177.49 \pm 4.69
UFGPOOL-SUM	77.77 \pm 2.60	81.59 \pm 1.40	80.92 \pm 1.68	77.88 \pm 1.24	78.80 \pm 0.56	41.74 \pm 0.84
UFGPOOL-SPECTRUM	77.23 \pm 2.40	82.05 \pm 1.28	79.83 \pm 1.88	77.54 \pm 2.24	78.36 \pm 0.77	41.67 \pm 0.95

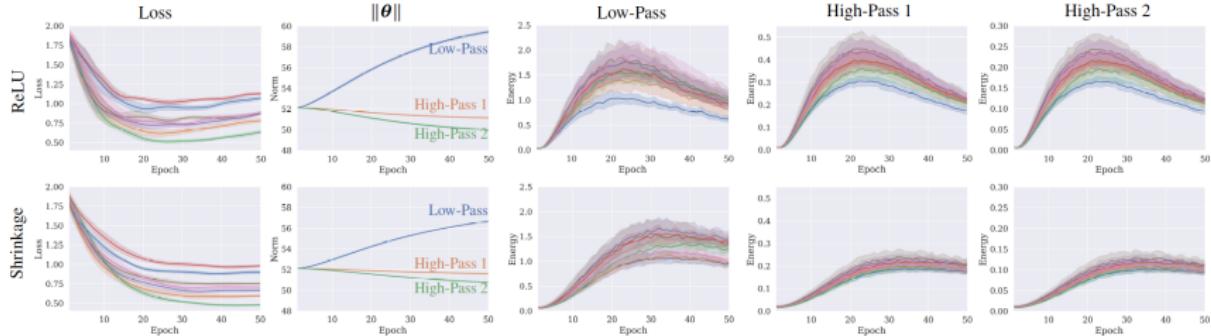
† The top three are highlighted by **First**, **Second**, **Third**.

Framelet spectrum, traning loss and network capacity



- Learning behavior of the final framelet convolutional layer of GNN with two UFG-CONV for Cora.
- The coefficients after shrinkage activation are proportional to the framelet power spectrum at the coefficient scale. We thus let the threshold level σ proportionate upon the framelet energy $\|\mathcal{W}_{r,j}X\|^2$ ($r > 0$) for high-passes. For example, the framelet spectrum curves in training show a higher magnitude order of the low-pass (column 3) than those of high-passes (columns 4-5). This is because coefficients in high-passes reflect more detailed characteristics than in low-pass.
- Compared with the ReLU case (row 1), the shrinkage activation (row 2) filters out some high-pass coefficients in graph convolution, which results in much smaller framelet spectra for high-passes. In contrast, the low-pass shrinkage involves no cutoff, and the energy is less distinguishable from the ReLU case.

Framelet spectrum, traning loss and network capacity



- The training loss curve of each output feature (column 1) indicates that shrinkage allows for **more stable** training, with a monotonically decreasing loss.
- The splitting in low-pass and high-passes for loss suggests a more flexible and precise control of the training. It also opens the possibility of designing a new **weighted loss** taking account of framelet spectrum.

Discussion

- We explore the adaptation of graph framelets for graph neural networks, and link graph neural networks and signal processing. In many node-level or graph-level tasks, framelet convolutions can reduce both feature and structure noises.
- We introduce shrinkage activation that thresholds high-pass coefficients in framelet convolution, which strengthens the network denoising capability and simultaneously compresses graph signal at a remarkable rate.
- The framelet pooling outperforms baselines on a variety of graph property prediction tasks.

Thank you!

References

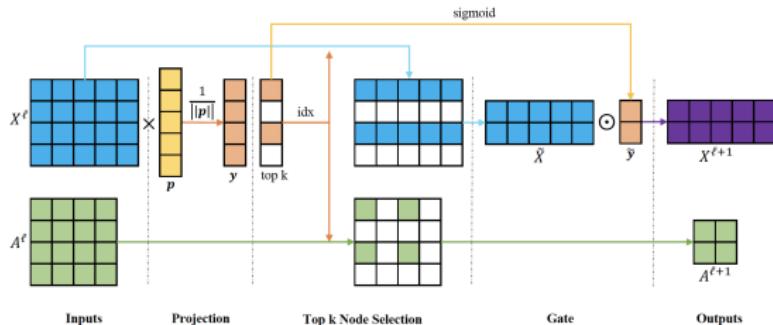
- How framelets Enhance Graph Neural Networks.
ICML 2021
- Decimated Framelet System on Graphs and Fast \mathcal{G} -Framelet Transforms. *JMLR 2021*
- Tight Framelets and Fast Framelet Filter Bank Transforms on Manifolds. *ACHA 2020*
- Haar Graph Pooling. *ICML 2020*
- Fast Haar Transforms for Graph Neural Networks.
Neural Networks 2020

Codes YuGuangWang@Github

Methods and Architecture in Comparison

- In the experiment, we normalize the label value by subtracting the mean and scaling the standard deviation (Std Dev) to 1. We then need to convert the predicted output to the original label domain (by re-scaling and adding the mean back).
- Following *Gilmer et al. (2017)*, we use mean squared error (MSE) as the loss for training and mean absolute error (MAE) as the evaluation metric for validation and test.
- Here, the splitting percentages for training, validation, and test are 80%, 10%, and 10%, respectively. We set the hidden dimension of the GCN layer as 64, the Adam for optimization with the learning rate 5.0e-4, and the maximal epoch 50 with no early stop. We do not use dropout as it would slightly lower the performance. For better comparison, we repeat all experiments ten times with different random seeds.
- We use the same GNN architecture to test HaarPool and SAGPool *Lee et al., (2019)*. : one GCN layer, one graph pooling layer, plus one 3-layer MLP.
- We compare the performance (test MAE) of the GCN-HaarPool against the GCN-SAGPool and other methods including Random Forest **RF** *Breiman, (2001)*, Multitask Networks **Multitask** *Ramsundar et al., (2015)*, Kernel Ridge Regression **KRR** *Cortes & Vapnik, (1995)*, Graph Convolutional models **GC** *Altae-Tran et al., (2017)*.

Existing Graph Pooling Methods



TopKPooling *Gao & Ji 2019, Knyazev et al 2019, Cangea, Liò et al 2018.*

- **Spatial Method**
 - Global (Non-hierarchical) pooling:
ChebNet, Set2Set, SortPool, MPNN
 - Hierarchical pooling:
DiffPool, TopKPool, SAGPool, EdgePool
- **Spectral Method**
 - LaPool, EigenPool, SOPool