

CS6208 : Advanced Topics in Artificial Intelligence

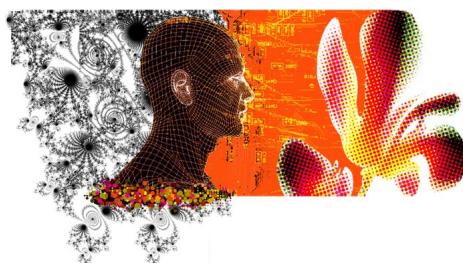
Graph Machine Learning

Lecture 8 : Benchmarking Graph Neural Networks and Graph Datasets

Semester 2 2022/23

Xavier Bresson

<https://twitter.com/xbresson>



Department of Computer Science
National University of Singapore (NUS)



Outline

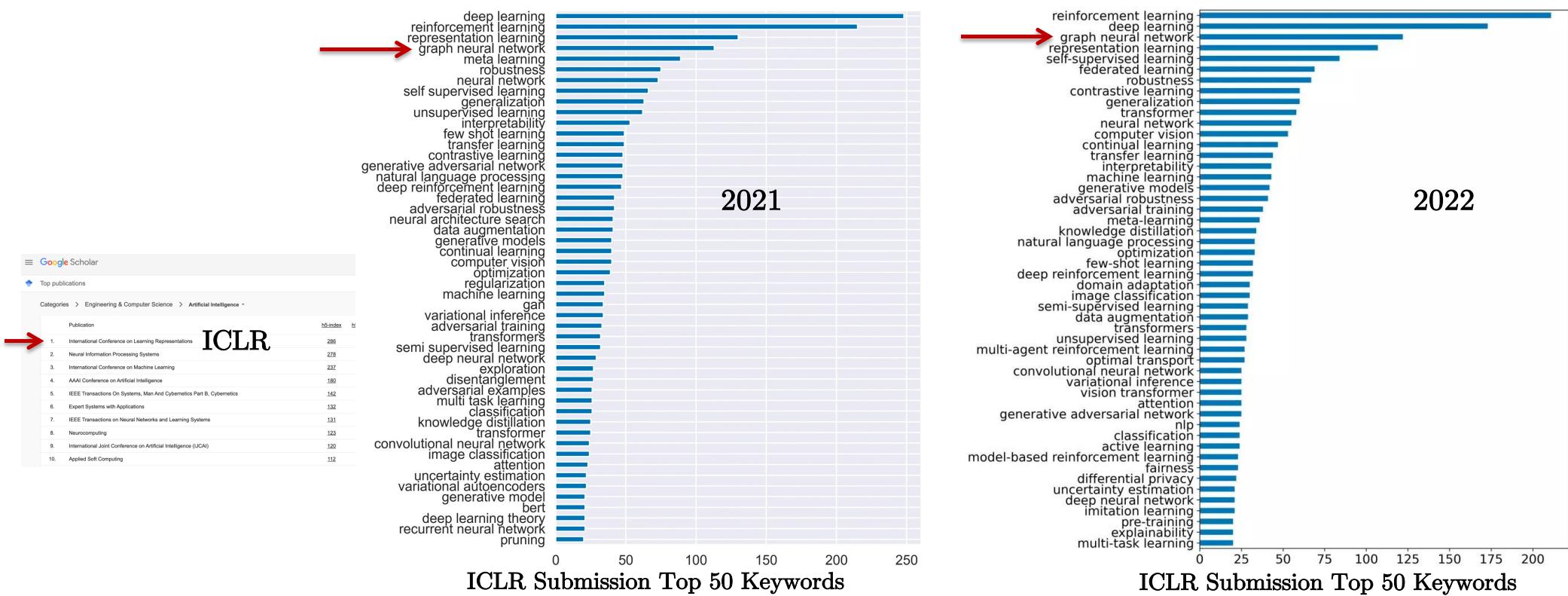
- Emerging field
- Tracking progress
- Benchmarking GNNs
 - Architectures
 - Datasets
 - Code infrastructure
 - Experimental setting
 - Results and insights
- Graph datasets
- Conclusion

Outline

- Emerging field
- Tracking progress
- Benchmarking GNNs
 - Architectures
 - Datasets
 - Code infrastructure
 - Experimental setting
 - Results and insights
- Graph datasets
- Conclusion

ICLR statistics

- GNNs have become the standard toolkit for analyzing and learning from data on graphs.
- One of the hottest machine learning topics since 2021.



Application domains

- GNNs are flexible to adapt to complex data structure and combine multi-modal features.
- Applicative domains :

- Chemistry^[1,2] (generate new drugs and materials)
- Physics^[3,4,5] (detect particles, accelerate physics)
- Recommender systems^[6,7] (leverage consumer-product choices)
- Social sciences^[8,9] (predict future friends, identify fake news)
- Knowledge graphs^[10,11] (reasoning with entity-relationship)
- Neuroscience^[12] (understanding brain mechanisms and neuro-degenerative diseases)
- Computer Vision^[13] (scene understanding for visual reasoning)
- Natural Language Processing^[14,15] (common sense reasoning)
- Combinatorial Optimization^[16,17] (better/faster approximated solutions to NP-hard problems)

[1] Duvenaud, Maclaurin, Iparraguirre, Bombarell, Hirzel, Aspuru-Guzik, Adams, Convolutional networks on graphs for learning molecular fingerprints, 2015

[2] Gilmer, Schoenholz, Riley, Vinyals, Dahl, Neural message passing for quantum chemistry, 2017

[3] Battaglia, Pascanu, Lai, Rezende, Daniil, Interaction networks for learning about objects, relations and physics, 2016

[4] Crammer, Xu, Battaglia, Ho, Learning symbolic physics with graph networks, 2019

[5] Sanchez-Gonzalez, Godwin, Pfaff, Ying, Leskovec, Battaglia, Learning to simulate complex physics with graph networks, 2020

[6] Monti, Bronstein, Bresson, Geometric matrix completion with recurrent multi-graph neural networks, 2017

[7] Ying, He, Chen, Eksombatchai, Hamilton, Leskovec, Graph convolutional neural networks for web-scale recommender systems, 2018

[8] Kipf, Welling, Semi-supervised classification with graph convolutional networks, 2017

[9] Monti, Frasca, Eynard, Mannion, Bronstein, Fake news detection on social media using geometric deep learning, 2019

[10] Schlichtkrull, Kipf, Bloem, VanDenBerg, Titov, Welling, Modeling relational data with graph convolutional networks, 2018

[11] Chami, Wolf, Juan Sala Ravi, Re, Low-dimensional hyperbolic knowledge graph embeddings, 2020

[12] Parisot, Ktena, Ferrante, Lee, Guerrero, Glocker, Rueckert, Disease prediction using graph networks: Application to Autism Spectrum Disorder and Alzheimer's disease, 2018

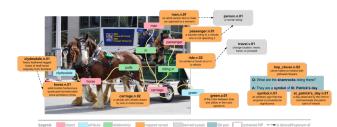
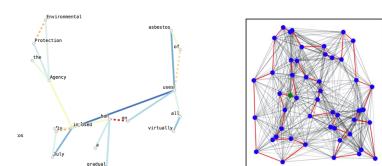
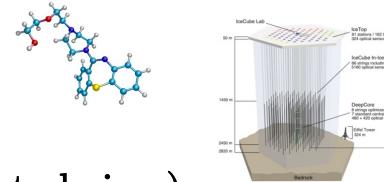
[13] Johnson, Gupta, Fei-Fei, Image generation from scene graphs, 2018

[14] Vaswani, Shazeer, Parmar, Uszkoreit, Jones, N Gomez, Kaiser, Polosukhin, Attention is all you need, 2017

[15] Yasunaga et-al, Deep Bidirectional Language-Knowledge Graph Pretraining, 2022

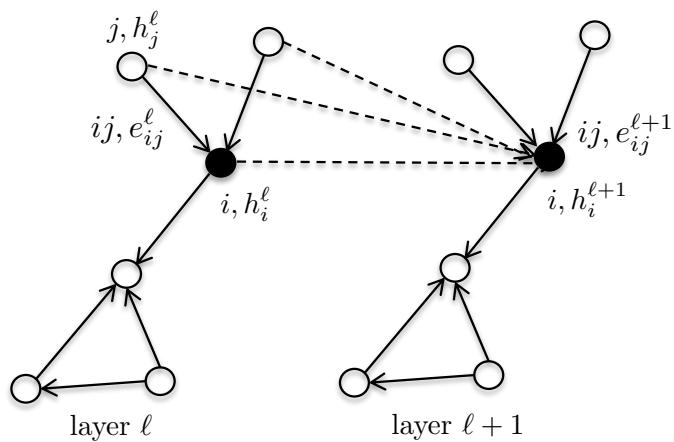
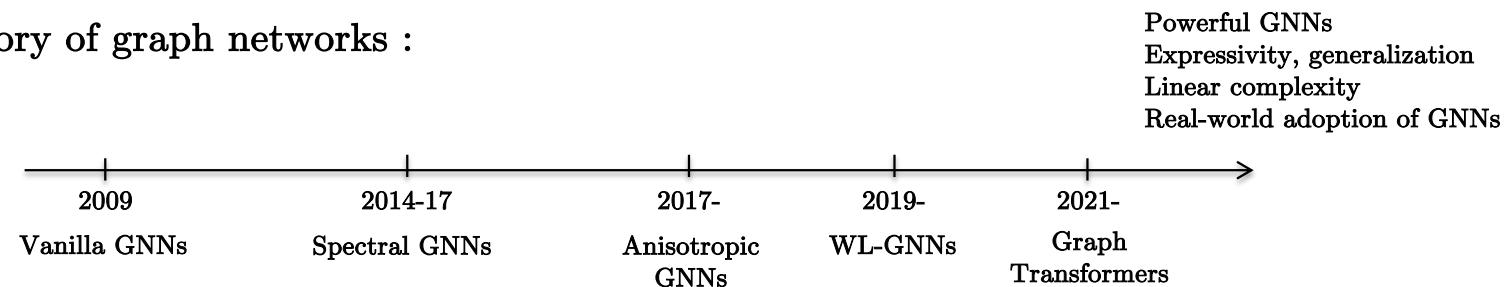
[16] Vinyals, Fortunato, Jaitly, Pionter, 2015

[17] Bello, Pham, Le, Norouzi, Bengio, Neural combinatorial optimization with reinforcement learning, 2016



A decade of GNN development

- Brief history of graph networks :



Node-update : $h_i^{\ell+1} = f_{\text{node}}(h_i^\ell, \{h_j^\ell, e_{ij}^\ell : j \in \mathcal{N}_i\}) \in \mathbb{R}^d$

Edge-update : $e_{ij}^{\ell+1} = f_{\text{edge}}(e_{ij}^\ell, h_i^\ell, h_j^\ell) \in \mathbb{R}^d$



"graph neural networks"



All Images Videos News

About 1,260,000 results (0.61 seconds)

Outline

- Emerging field
- Tracking progress
- Benchmarking GNNs
 - Architectures
 - Datasets
 - Code infrastructure
 - Experimental setting
 - Results and insights
- Graph datasets
- Conclusion

Progress or not

- The field has grown extensively in the last few years 😊
- Unfortunately, evaluating the effectiveness of new architectures/ideas has become difficult for two major reasons.
 - We have been evaluated progress on small datasets such as Cora^[1], Citeseer^[2] and TU^[3].
 - Cora is a single graph of 2.7K nodes, TU-IMDB has 1.5K graphs with 13 nodes on average, and TU-MUTAG has 188 molecules with 18 nodes.
 - Simple or graph-agnostic architectures provide statistically same performance as complex architectures^[4,5] (74.2 ± 6.3 vs. 75.6 ± 7.2).
 - We have not rigorously enforced standardized experimental settings for fair comparisons between models^[6].
- It has become critical to solve these issues to trust the current models and new ideas.

[1] McCallum, Nigam, Rennie, Seymore, Automating the construction of internet portals with machine learning, 2000

[2] Getoor, Link-based classification, 2005

[3] Kersting, Kriege, Morris, Mutzel, Neumann, Benchmark data sets for graph kernels, 2020

[4] Hoang, Maehara, Revisiting graph neural networks, 2019

[5] Chen, Bian, Sun, Are powerful graph neural nets necessary? a dissection on graph classification, 2019

[6] Errica, Podda, Bacciu, Micheli, A fair comparison of graph neural networks for graph classification, 2019

Benchmarking

- How to identify and quantify what types of architectures, first principles are universal, generalizable, and scalable to larger and more challenging datasets?
- Theoretical approach : Assumptions on data, design model, prove expressivity and generalization bound (latter is hard because it requires to simultaneously analyze data property, architecture, loss energy landscape and optimization).
- Empirical approach : Benchmarking has been beneficial for driving progress, identifying essential ideas, and solving domain-specific problems^[1].
 - The 2012 ImageNet challenge^[2] has provided a benchmark dataset that has triggered the deep learning revolution.
- Challenges
 - Developing a rigorous experimental setting for fair comparisons and being reproducible.
 - Using datasets that can statistically separate model performance.
 - Benchmarking distinct fundamental graph tasks.

[1] Weber, Saelens, Cannoodt, Soneson, Hapfelmeier, Gardner, Boulesteix, Saeys, Robinson, Essential guidelines for computational method benchmarking, 2019

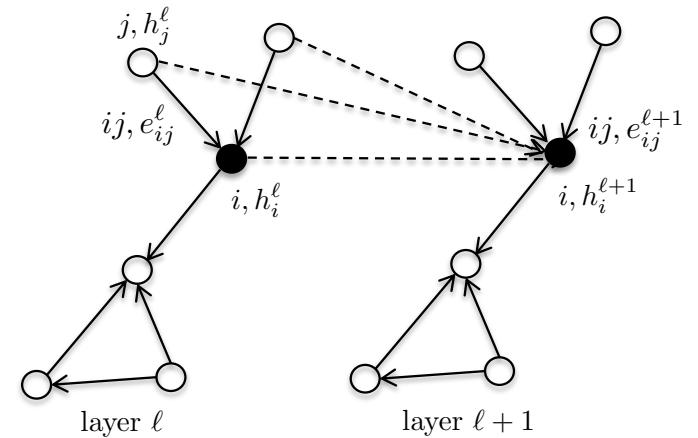
[2] Deng, Dong, Socher, Li, Li, Fei-Fei, ImageNet: A Large-Scale Hierarchical Image Database, 2009

Outline

- Emerging field
- Tracking progress
- Benchmarking GNNs
 - Architectures
 - Datasets
 - Code infrastructure
 - Experimental setting
 - Results and insights
- Graph datasets
- Conclusion

Message Passing GNNs

- The most popular and oldest class of GNNs – obsolete?
- What are the minimal inner structures to design a MP-GNN?
 - Invariant by node permutation (equivariant/invariant layers)
 - Independent of graph size n and neighborhood size
 - Locality (local reception field - only neighbors are considered)
 - Graph convolution operator (weight sharing across graph)
 - Linear complexity $O(E)$, E being the number of edges (reducing to $O(n)$ for sparse/real graphs)



Node/edge-update :

$$h_i^{\ell+1} = f_{\text{node}}(h_i^\ell, \{h_j^\ell, e_{ij}^\ell : j \in \mathcal{N}_i\}) \in \mathbb{R}^d$$

$$e_{ij}^{\ell+1} = f_{\text{edge}}(e_{ij}^\ell, h_i^\ell, h_j^\ell) \in \mathbb{R}^d$$

Isotropic MP-GNNs

- We consider the class of isotropic MP-GNNs when the node update equation treats every “edge direction” equally, i.e. each neighbor contributes equally to the update of the central node :

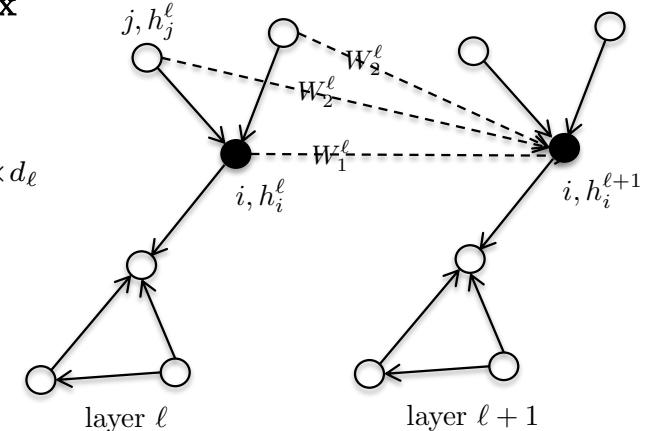
$$h_i^{\ell+1} = \sigma \left(W_1^\ell h_i^\ell + \sum_{j \in \mathcal{N}_i} W_2^\ell h_j^\ell \right),$$

ReLU Sum/mean/max

$$h^{\ell+1} \in \mathbb{R}^{n \times d_{\ell+1}}, \quad h^\ell \in \mathbb{R}^{n \times d_\ell}, \quad W_{1,2}^\ell \in \mathbb{R}^{d_{\ell+1} \times d_\ell}$$

- Models :

- GCNs^[1,2]
- GraphSage^[3]
- ChebNets^[4]



- Isotropic MP-GNNs are limited to learn two template weights; one for the central node h_i (weight W_1) and one for the 1-hop neighbors h_j (weight W_2).

[1] Sukhbaatar, Szlam, Fergus, Learning multiagent communication with backpropagation, 2016

[2] Kipf, Welling, Semi-supervised classification with graph convolutional networks, 2017

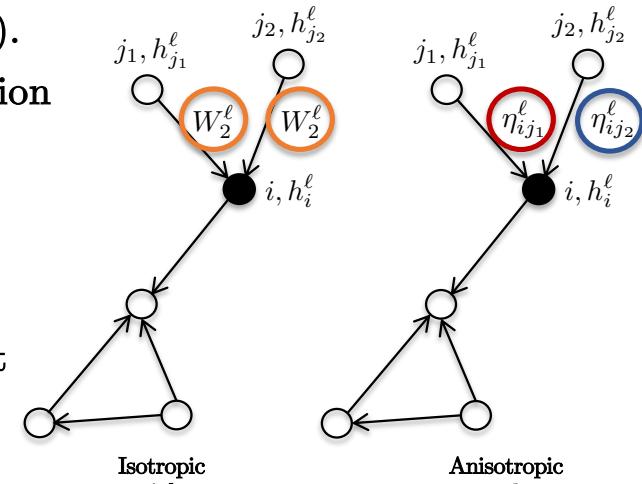
[3] Hamilton, Ying, Leskovec, Inductive representation learning on large graphs, 2017

[4] Defferrard, Bresson, Vandergheynst, Convolutional neural networks on graphs with fast localized spectral filtering, 2016

Anisotropic MP-GNNs



- MP-GNNs such as vanilla GCNs, GraphSage, ChebNets compute isotropic filters as there is no notion of direction on arbitrary graphs.
- How to get anisotropy back?
 - Natural edge features^[1] if available (e.g. different bond connections between atoms in molecular graphs or distinct edge attributes in knowledge graphs).
 - Learn anisotropy with a mechanism invariant by index permutation
 - MoNets^[2] with edge degrees
 - GatedGCNs^[3] with edge gates
 - GAT^[4] with attention mechanism^[7]
- We instantiate an anisotropic MP-GNN with an update equation that treats every edge direction differently



Learnable anisotropic weights

$$h_i^{\ell+1} = \sigma \left(W_1^\ell h_i^\ell + \sum_{j \in \mathcal{N}_i} \eta_{ij}^\ell W_2^\ell h_j^\ell \right), \quad h^{\ell+1} \in \mathbb{R}^{n \times d_{\ell+1}}, \quad h^\ell \in \mathbb{R}^{n \times d_\ell}, \quad W_{1,2}^\ell \in \mathbb{R}^{d_{\ell+1} \times d_\ell},$$

[1] Gilmer, Schoenholz, Riley, Vinyals, Dahl, Neural message passing for quantum chemistry, 2017

[2] Monti, Bosacini, Masci, Rodolà, J. Svoboda, M. Bronstein, Geometric deep learning on graphs and manifolds using mixture model CNNs, 2016

[3] Bresson, Laurent, Residual gated graph convnets, 2017

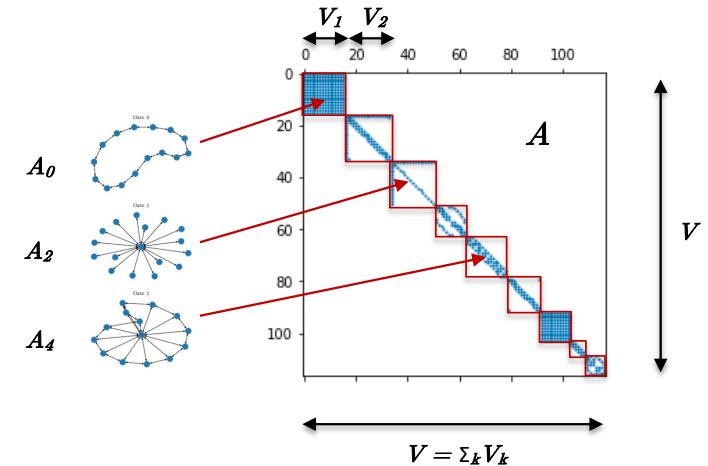
[4] Velickovic, Cucurull, Casanova, Romero, Lio, Bengio, Graph Attention Networks, 2018

Batch normalization and residual connection

- MP-GNNs benefit from BN^[1,2] and RC^[1,3]
 - Speed up learning process
 - Improve performance with better generalization
- How to batch graphs of different sizes ?
 - For images of size $B \times V_x \times V_y \times d$, BN is performed along the batch dimension, i.e. dim=0.
 - For graphs, a (big) sparse block diagonal matrix A of size $V \times d$ is first built from K matrices A_k and BN is carried out along the node direction, i.e. dim=0.

$$h_i^{\ell+1} = h_i^\ell + \sigma\left(\text{BN}\left(\hat{h}_i^{\ell+1}\right)\right)$$

$$\hat{h}_i^{\ell+1} = f_{\text{MP-GNN}}\left(h_i^\ell, \{h_j^\ell : j \in \mathcal{N}_i\}\right)$$



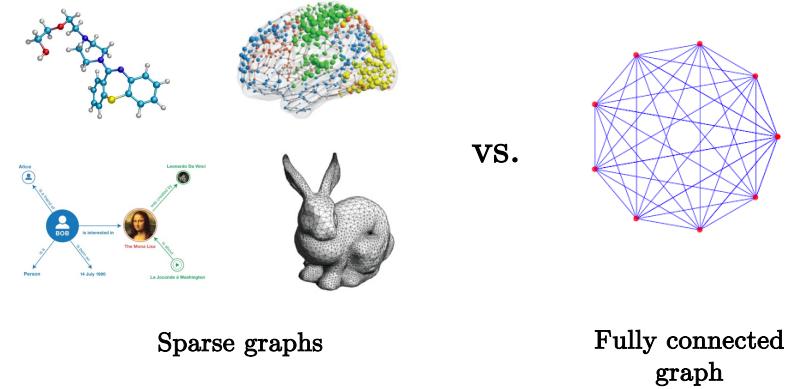
[1] Bresson, Laurent, Residual gated graph convnets, 2017

[2] Ioffe, Szegedy, Batch normalization: Accelerating deep network training by reducing internal covariate shift, 2015

[3] Li, Muller, Qian, Delgadillo, Abualshour, Thabet, Ghanem, Deepgcns: Making gcns go as deep as cnns, 2019

Sparsity and local computations

- MP-GNNs leverage graph sparsity
 - Sparsity is a great inductive bias for generalization
 - Sparse vs full graphs
- Local computations
 - Node/edge update equations are local, only depending on neighborhood \mathcal{N}_i of node i , and independent of graph size n , making the space/time complexity $O(n)$ for sparse graphs.
 - MP-GNNs are naturally parallelizable, and fast sparse matrix multiplications are implemented via GPU-based libraries s.a. DGL^[1], PyG^[2], TensorFlowGNNs^[3].



Sparse graphs

Fully connected graph

$$h_i^{\ell+1} = f_{\text{node}}(h_i^\ell, \{h_j^\ell, e_{ij}^\ell : j \in \mathcal{N}_i\}) \in \mathbb{R}^d$$
$$e_{ij}^{\ell+1} = f_{\text{edge}}(e_{ij}^\ell, h_i^\ell, h_j^\ell) \in \mathbb{R}^d$$

DGL  **PyTorch geometric**

 **TensorFlow GNN**

[1] Wang-etal, Deep graph library: Towards efficient and scalable deep learning on graphs, 2019

[2] Fey, Lenssen, Fast graph representation learning with pytorch geometric, 2019

[3] Oleksandr et-al, TF-GNN: Graph Neural Networks in TensorFlow, 2022

Graph isomorphism networks

- GINs^[1] : MP-GNN as expressive as the 1-WL test to distinguish non-isomorphic graphs.
- Node aggregation of the form :

$$f_{\text{NN}}(h_i^\ell, \{h_j^\ell\}_{j \in \mathcal{N}_i}) = (1 + \varepsilon)g(h_i^\ell) + \sum_{j \in \mathcal{N}_i} g(h_j^\ell)$$

irrational injective sum

- It is difficult to design an analytical injective function \Rightarrow MLP is used to approximate g (existence guaranteed by the universal approximation theorem) :

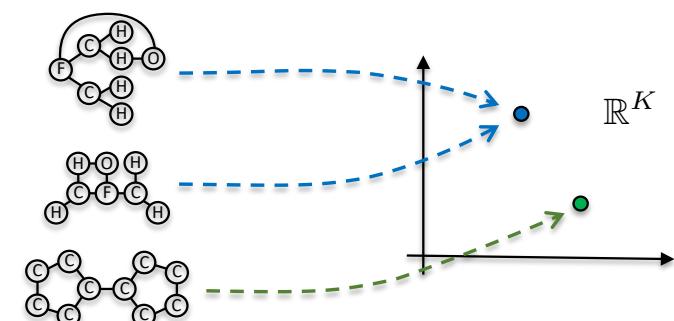
$$h_i^{\ell+1} = f_{\text{GIN}}(h_i^\ell, \{h_j^\ell\}_{j \in \mathcal{N}_i}) = \text{MLP}^\ell \left((1 + \varepsilon)h_i^\ell + \sum_{j \in \mathcal{N}_i} h_j^\ell \right)$$

- Graph readout function must also be injective :

$$h_G = \text{MLP} \left(\sum_{i \in V} h_i^L \right) \in \mathbb{R}^K$$

sum

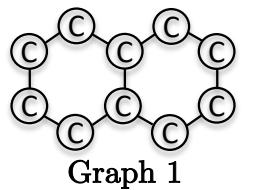
- Pioneer work on GNN expressivity



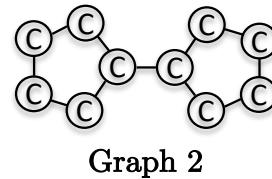
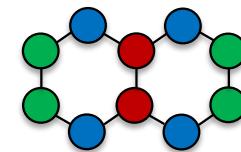
[1] Xu, Hu, Leskovec, Jegelka, How powerful are graph neural networks?, 2019

Limitation of 1-WL test

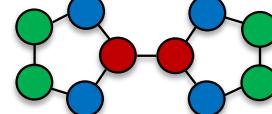
- WL test^[1] is not a sufficient condition, it can fail to differentiate simple non-isomorphic graphs.
 - Example of two simple distinct graphs with same color signatures but not isomorphic



⇒



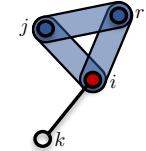
⇒



- Need to improve the expressivity of the original WL test

[1] Weisfeiler, Lehman, A reduction of a graph to a canonical form and an algebra arising during this reduction, 1968

Equivariant GNNs

- How to design GNNs with the same expressivity power as the k-WL test?
 - Let us define a k-order equivariant GNNs^[1] :
- 

$y = m_{W^{L+1}} \circ g_{W^L} \circ \sigma \circ f_{W^{L-1}} \circ \sigma \circ f_{W^{L-2}} \circ \dots \sigma \circ f_{W^0} \circ h_0$
 with $k = \max_{\ell \in [0, L-1]} k_\ell$
 and $f_{W^\ell} : \mathbb{R}^{n^{k_\ell} \times d_\ell} \rightarrow \mathbb{R}^{n^{k_\ell+1} \times d_{\ell+1}}$

$$f_W(P \circ h) = P \circ f_W(h)$$

Equivariant linear layers

$$g_W(P \circ h) = g_W(h) \in \mathbb{R}^K, K \geq 1$$

Invariant linear layer
- Theorem : There exist k-order E-GNNs that can distinguish non-isomorphic graphs with the k-WL test.
 - However, k-order E-GNNs require $O(n^k)$ memory/speed complexities.
 - Note that we need at least $k=3$, therefore $O(n^3)$, to be more powerful than GINs as GINs have the discriminative power of the 1-WL/2-WL test^[2].

[1] Maron, Ben-Hamu, Shamir, Lipman, Invariant and equivariant graph networks, 2019

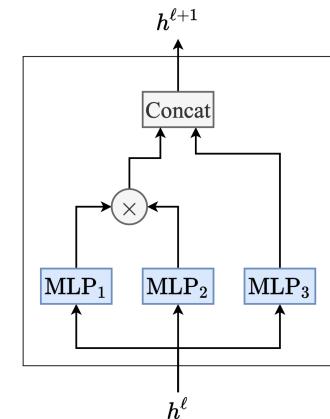
[2] Cai, Furer, Immerman, An optimal lower bound on the number of variables for graph identification, 1992

3-WL GNNs

- How to design GNNs that are 3-WL expressive but do not require $O(n^3)$ memory complexity?
- 3-WL GNNs^[1] : To achieve higher interactions between nodes, it is sufficient to multiply second-order tensors feature-wise.
 - Theorem : There exist 3-WL GNNs as expressive as the 3-WL test.
 - Memory is quadratic $O(n^2)$ but matrix multiplication implies $O(n^3)$ speed.
 - Also, matrix multiplication densifies sparse matrix.
 - This is one of the simplest, most scalable and most expressive GNNs in terms of 3-WL test.
- 3-WL GNN update layer :

$$h^{\ell+1} = \text{Concat}(m_{W_1^\ell}(h^\ell) \cdot m_{W_2^\ell}(h^\ell), m_{W_3^\ell}(h^\ell))$$

where $h^{\ell+1} \in \mathbb{R}^{n \times n \times d_{\ell+1}}$, $h^\ell \in \mathbb{R}^{n \times n \times d_\ell}$, $W_1, W_2, W_3 \in \mathbb{R}^{2 \times d_\ell \times d_{\ell+1}}$



[1] Maron, Ben-Hamu, Serviansky, Lipman, Provably powerful graph networks, 2019

RingGNNs

- Related method to 3-WL GNNs to design more expressive GNNs than GINs.
 - Higher-order interaction between nodes is produced by multiplying equivariant linear layers^[2].
 - RingGNN^[1] update layer :

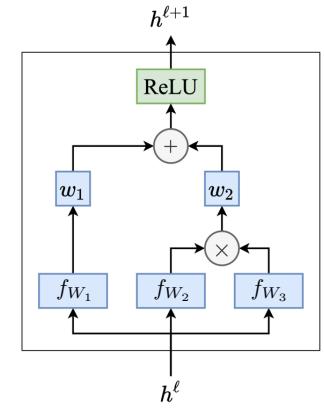
$$h^{\ell+1} = \sigma(w_1^\ell f_{W_1^\ell}(h^\ell) + w_2^\ell f_{W_2^\ell}(h^\ell) \cdot f_{W_3^\ell}(h^\ell)),$$

where $h^{\ell+1} \in \mathbb{R}^{n \times n \times d_{\ell+1}}$, $h^\ell \in \mathbb{R}^{n \times n \times d_\ell}$, $w_{1,2}^\ell \in \mathbb{R}$, $W_1, W_2, W_3 \in \mathbb{R}^{d_\ell \times d_{\ell+1} \times 17}$,

where f_W are the linear equivariant layers defined as :

$$(f_W(h))_{\cdot,\cdot,q'} = \sum_{k=1}^{15+2} \sum_{q=1}^{d_\ell} W_{k,q,q'} f_k(h_{\cdot,\cdot,q}) \in \mathbb{R}^{n \times n \times d_{\ell+1}}$$

and f_k are all possible 15 equivariant linear functions $f_k : \mathbb{R}^{n \times n} \rightarrow \mathbb{R}^{n \times n}$ for a given tensor $h \in \mathbb{R}^{n \times n}$, and 2 bias functions.



[1] Chen, Villar, Chen, Bruna, On the equivalence between graph isomorphism testing and function approximation with gnn, 2019

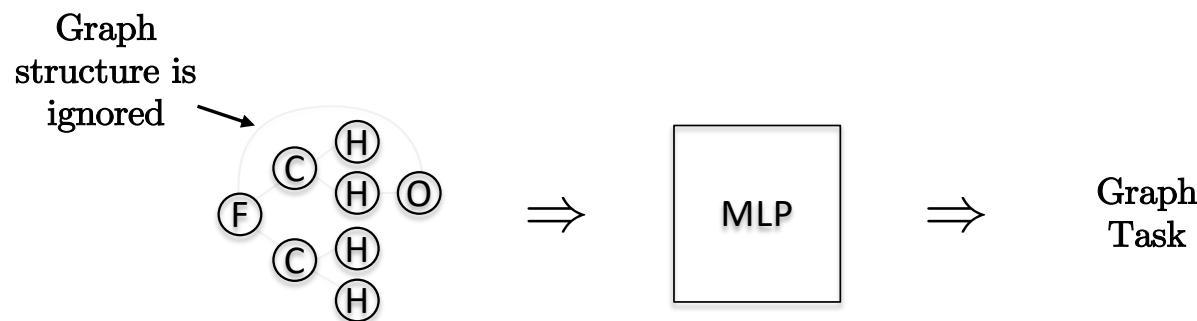
[2] Maron, Ben-Hamu, Shamir, Lipman, Invariant and equivariant graph networks, 2019

MLP baseline

- Graph-agnostic NNs : As a sanity check, we compare GNNs to a simple MLP network which updates each node independent of one-other :

$$h_i^{\ell+1} = \sigma(W^\ell h_i^\ell), \quad h^{\ell+1} \in \mathbb{R}^{n \times d_{\ell+1}}, \quad h^\ell \in \mathbb{R}^{n \times d_\ell}, \quad W^\ell \in \mathbb{R}^{d_{\ell+1} \times d_\ell}$$

and passes these features to the task-based readout layer.



Outline

- Emerging field
- Tracking progress
- Benchmarking GNNs
 - Architectures
 - Datasets
 - Code infrastructure
 - Experimental setting
 - Results and insights
- Graph datasets
- Conclusion

Designing datasets

- The main issue with datasets used in most GNN papers is their small size.
- What are the best properties for a dataset?
 - Representative, realistic, and medium/large-scale size
- Challenges
 - What theoretical tool to define the quality of dataset or validate its statistical representativeness for a given task?
 - What are the best features? Several arbitrary choices when preparing graphs, such as node and edge features. For example, e-commerce product features.
 - What size? Appropriate size may depend on the task complexity as well as the dimensionality and statistics of underlying data. Besides, we wish to use the smallest dataset that is representative of the task because GPUs are expensive and scarce resources.
 - How to quantify and identify noise/error in real datasets? Synthetic datasets control the noise but are not real. Real datasets are noisy.
- In practice, there is no perfect dataset.

Proposed datasets

- Appropriate datasets :
 - Datasets that are able to statistically separate the performance of GNNs (no more 74.2 ± 6.3 vs. 75.6 ± 7.2).
- Summary of the 7 medium-scale datasets used in the benchmark^[1] :

Domain & Construction	Dataset	#Graphs	#Nodes	Total #Nodes	Task
Chemistry: Real-world molecular graphs	ZINC	12K	9-37	277,864	Graph Regression
Mathematical Modelling: Artificial graphs generated from Stochastic Block Models	PATTERN	14K	44-188	1,664,491	Node Classification
	CLUSTER	12K	41-190	1,406,436	
Computer Vision: Graphs constructed with SLIC super-pixels of images	MNIST	70K	40-75	4,939,668	Graph Classification
	CIFAR10	60K	85-150	7,058,005	
Combinatorial Optimization: Uniformly generated artificial Euclidean graphs	TSP	12K	50-500	3,309,140	Edge Classification
Social Networks: Real-world citation graph	COLLAB	1	235,868	235,868	Edge Classification
Circular Skip Links: Isomorphic graphs with same degree	CSL	150	41	6,150	Graph Classification

4 datasets are artificially generated, 2 datasets are semi-artificial, and 2 are real-world.

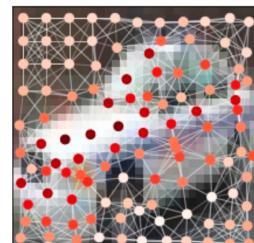
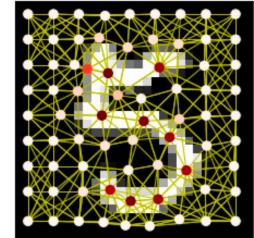
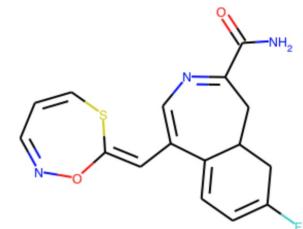
Sizes in terms of total number of nodes vary between 0.27M to 7M.

4 most fundamental graph tasks : graph regression, graph classification, node classification and link prediction.

[1] Dwivedi, Joshi, Laurent, Bengio, Bresson, Benchmarking graph neural networks, 2020

Proposed datasets

- ZINC^[1] : A popular real-world molecular dataset of 250K graphs, out of which we randomly select 12K for efficiency. We consider the task of graph property regression for constrained solubility, an important chemical property for designing generative GNNs for molecules^[2].
 - Statistics : 10,000 train/1,000 validation/1,000 test graphs of sizes 9-37 nodes/heavy atoms.
- MNIST^[3]/CIFAR10^[4] : Classical image classification datasets converted into graphs using super-pixels^[5] and assigning node features as the super-pixel coordinates and mean intensity. These datasets are sanity-checks, as we expect most GNNs to perform close to 100% for MNIST and well enough for CIFAR10.
 - Statistics : MNIST has 55,000 train/5,000 validation/10,000 test graphs of sizes 40-75 nodes and CIFAR10 has 45,000 train/5,000 validation/10,000 test graphs of sizes 85-150 nodes.



[1] Irwin, Sterling, Mysinger, Bolstad, Coleman, Zinc: a free tool to discover chemistry for biology, 2012

[2] Jia, Lin, Ying, You, Leskovec, Aiken, Redundancy-free computation graphs for graph neural networks, 2019

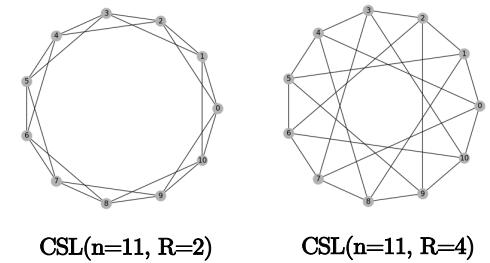
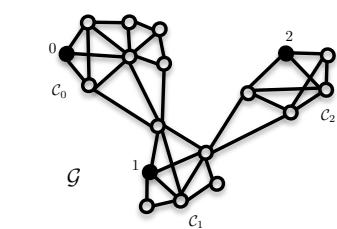
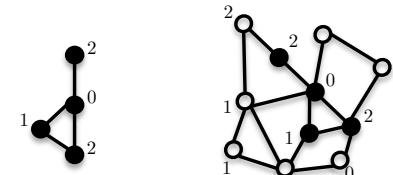
[3] LeCun, Bottou, Bengio, Haffner, Gradient-based learning applied to document recognition, 1998

[4] Alex Krizhevsky et al, Learning multiple layers of features from tiny images, 2009

[5] Achanta, Shaji, Smith, Lucchi, Fua, Süsstrunk, Slic superpixels compared to state-of-the-art superpixel methods, 2012

Proposed datasets

- **PATTERN/CLUSTER** : Node classification tasks generated with Stochastic Block Models^[1], which are widely used to model communities in social networks by modulating the intra- and extra-communities connections. PATTERN tests the fundamental task of recognizing specific predetermined subgraphs^[2].
 - Statistics : PATTERN has 10,000 train/2,000 validation/2,000 test graphs of sizes 50-180 nodes. CLUSTER has 10,000 train/1,000 validation/1,000 test graphs of sizes 40-190 nodes.
- **CSL**^[3] : Synthetic to test the expressivity of GNNs. Graphs are isomorphic if they have the same degree and the task is to classify non-isomorphic graphs.
 - Statistics : 5-fold cross-validation split of 150 graphs of sizes 41 nodes with train-val-test ratio 3:1:1.



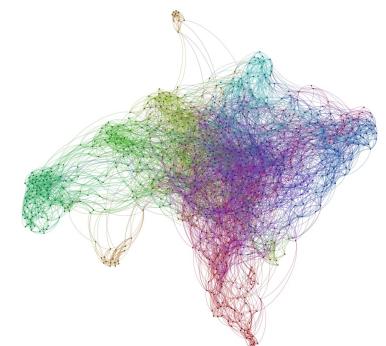
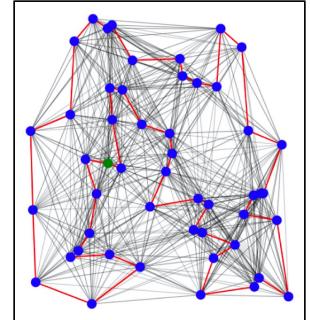
[1] Abbe, Community detection and stochastic block models: recent developments, 2017

[2] Scarselli, Gori, Tsoi, Hagenbuchner, Monfardini, The Graph Neural Network Model, 2009

[3] Murphy, Srinivasan, Rao, Ribeiro, Relational pooling for graph representations, 2019

Proposed datasets

- Traveling Salesman Problem (TSP) : Link prediction on 2D Euclidean graphs to identify edges belonging to the optimal TSP solution given by Concorde^[1]. TSP is the most studied NP-hard combinatorial problem with a growing body of literature on leveraging GNNs to learn better solvers^[2,3].
 - Statistics : 10,000 train TSPs / 1,000 validation TSPs / 1,000 test TSPs, where the number of nodes is randomly selected in [50, 500].
- OGB-COLLAB^[4] : Link prediction dataset proposed by OGB corresponding to a collaboration network between scientists. The task is to predict future author collaboration relationships given past collaboration links.
 - Statistics : A single large temporal graph of size 235K nodes with given train/validation/test edge splits.



[1] Applegate, Bixby, Chvatal, Cook, Concorde tsp solver, 2006

[2] Vinyals, Fortunato, Jaity, Pointer, 2015

[3] Bello, Pham, Le, Norouzi, Bengio, Neural combinatorial optimization with reinforcement learning, 2016

[4] Hu, Fey, Zitnik, Dong, Ren, Liu, Catasta, Leskovec, Open graph benchmark: Datasets for machine learning on graphs, 2020

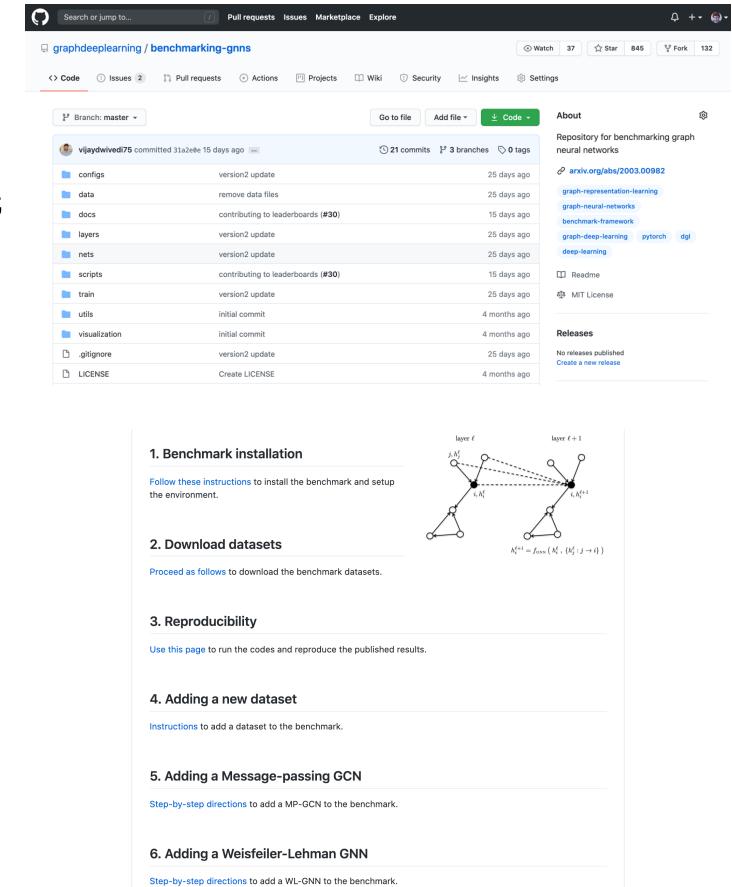
Outline

- Emerging field
- Tracking progress
- Benchmarking GNNs
 - Architectures
 - Datasets
 - **Code infrastructure**
 - Experimental setting
 - Results and insights
- Graph datasets
- Conclusion

Benchmark infrastructure

PyTorch DGL

- GitHub Repo (2.1k+ stars, 4k+ forks as of Mar 2023)
<https://github.com/graphdeeplearning/benchmarking-gnns>
- Objectives
 - Ease-of-use and modular, enabling new users to experiment and study the building blocks of GNNs.
 - Experimentally rigorous and fair for all models being benchmarked.
 - Being comprehensive for tracking progress of new GNNs and novel dataset/task.
- Components
 - Data pipeline
 - GNN layers and models
 - Training and evaluation functions
 - Network and hyperparameter configurations
 - Scripts for reproducibility



The screenshot shows the GitHub repository page for `graphdeeplearning/benchmarking-gnns`. The repository has 37 stars and 132 forks. The code tab is selected, showing a list of commits from `vijayvwivedi75` over the past 15 days. The repository is described as a "Repository for benchmarking graph neural networks" with tags for arXiv.org, PyTorch, and DGL. Below the repository details, there is a "Benchmarking Guide" with six steps: 1. Benchmark installation, 2. Download datasets, 3. Reproducibility, 4. Adding a new dataset, 5. Adding a Message-passing GCN, and 6. Adding a Weisfeiler-Lehman GNN. Each step includes instructions and diagrams illustrating the process. A diagram at the top right shows two layers of nodes, \mathcal{N}^ℓ and $\mathcal{N}^{\ell+1}$, with arrows indicating the flow of information between them.

Outline

- Emerging field
- Tracking progress
- Benchmarking GNNs
 - Architectures
 - Datasets
 - Code infrastructure
 - **Experimental setting**
 - Results and insights
- Graph datasets
- Conclusion

Experimental setting

- Data splits : Given for ZINC, MNIST, CIFAR10, CSL, OGB-COLLAB, random for PATTERN, CLUSTER, TSP.
- Training :
 - Adam optimizer w/ learning rate decay strategy.
Initial learning rate is $1e-3/1e-4$, reduced by half if validation loss does not improve after 5/10 epochs.
Training is stopped if learning rate $\leq 1e-6$ or computational time ≥ 12 hours.
 - Statistics with 4 results using 4 different seeds are reported.
- Parameter budgets :
 - Our goal is not to find the optimal hyperparameters (computationally expensive) but to compare models within the same parameter budget.
 - Two parameter budgets :
 - 100k (w/ $L=4$) and 500k (w/ $L=16$) parameters for each GNNs for all tasks.

Outline

- Emerging field
- Tracking progress
- Benchmarking GNNs
 - Architectures
 - Datasets
 - Code infrastructure
 - Experimental setting
 - Results and insights
- Graph datasets
- Conclusion

Benchmarking results

- Result #1 : MP-GNNs outperformed WL-GNNs on all datasets.
 - Potential reasons :
 - $O(n^2)/O(n^3)$ memory/speed complexities of WL-GNNs do not scale to medium-scale datasets.
 - Best result for ZINC, which is the dataset with the smallest sizes, $n \in [9,37]$.
 - Early stage of developments of WL-GNNs.
- Result #2 : MP-GNNs benefit from batch normalization and residual connection :
 - Boost performance
 - Accelerate training
 - WL-GNNs do not benefit from RC, BN or LN (performance degrades).

Model	L	#Param	Test Acc. \pm s.d.	PATTERN		#Epoch	NODE CLASSIFICATION		#Param	Test Acc. \pm s.d.	#Epoch	CLUSTER
				Train Acc. \pm s.d.	Epoch		Epoch/Total	#Param	Test Acc. \pm s.d.	Epoch	Epoch/Total	
MLP	4	10528	85.48 \pm 0.02	85.19 \pm 0.07	85.17 \pm 0.13	42	8.95 \pm 0.01	106455	53.44 \pm 0.05	20.93 \pm 0.02	42.25	8.57 \pm 0.01
GCN	4	100023	63.80 \pm 0.04	65.76 \pm 0.13	65.00 \pm 0.09	11000	11.00 \pm 0.51hr	501657	53.44 \pm 0.05	34.01 \pm 0.02	79.75	65.72 \pm 0.08hr
GAT	4	100023	71.89 \pm 0.34	78.49 \pm 1.92	81.50	492.19 \pm 1.13hr	501687	68.49 \pm 0.076	71.72 \pm 2.212	79.75	77.09 \pm 0.08hr	
GraphSage	4	101739	59.16 \pm 0.006	50.47 \pm 0.014	43.75	93.41 \pm 1.17hr	102187	50.45 \pm 0.145	54.37 \pm 0.203	64.00	53.56 \pm 0.97hr	
	16	502842	50.49 \pm 0.001	50.48 \pm 0.001	46.50	391.16 \pm 1.19hr	303350	63.84 \pm 0.110	86.71 \pm 0.167	57.75	225.61 \pm 0.70hr	
MoNet	4	103775	85.48 \pm 0.037	85.50 \pm 0.044	89.75	35.71 \pm 0.90hr	104227	58.96 \pm 0.131	58.45 \pm 0.183	76.25	24.04 \pm 0.52hr	
GAT	4	109936	78.24 \pm 1.823	77.88 \pm 1.632	96.00	20.92 \pm 0.57hr	110700	57.73 \pm 0.323	58.33 \pm 0.342	67.25	19.17 \pm 0.27hr	
GatedGCN	4	1026990	78.71 \pm 0.161	90.21 \pm 0.476	53.50	50.33 \pm 0.77hr	527874	70.58 \pm 0.047	76.07 \pm 1.362	73.50	35.94 \pm 0.53hr	
	16	502723	85.58 \pm 0.088	86.00 \pm 0.073	65.25	644.71 \pm 1.91hr	502615	73.84 \pm 0.126	87.88 \pm 0.908	60.00	400.07 \pm 6.81hr	
GatedGCN-PE	16	502457	86.59 \pm 0.085	86.80 \pm 0.133	65.75	647.94 \pm 1.91hr	304253	76.08 \pm 0.136	88.91 \pm 0.720	57.75	399.66 \pm 6.58hr	
GIN	4	100884	85.59 \pm 0.011	85.85 \pm 0.030	93.00	15.24 \pm 0.40hr	103544	58.38 \pm 0.236	59.48 \pm 0.337	74.75	10.71 \pm 0.33hr	
RingGNN	2	504766	86.24 \pm 0.025	86.10 \pm 0.021	72.00	95.97 \pm 1.21hr	524202	22.34 \pm 0.000	22.30 \pm 0.000	43.25	501.84 \pm 6.22hr	
3WLGNN	3	100449	85.72 \pm 0.153	84.89 \pm 0.337	104.75	30.63 \pm 0.10hr	514380	Diverged	Diverged	Diverged	Diverged	
	8	502872	85.34 \pm 0.207	85.27 \pm 0.198	81.75	42.42 \pm 0.56hr	507252	53.77 \pm 0.139	87.88 \pm 0.907	51.00	210.52 \pm 5.22hr	
	8	581716	Diverged	Diverged	Diverged	586788	Diverged	Diverged	Diverged	Diverged	519.98 \pm 5.70hr	
GRAPH CLASSIFICATION												
Model	L	#Param	Test Acc. \pm s.d.	MINIST		#Epoch	CIFAR10		#Param	Test Acc. \pm s.d.	#Epoch	Epoch/Total
				Train Acc. \pm s.d.	Epoch		Epoch/Total	#Param	Test Acc. \pm s.d.	Epoch	Epoch/Total	
MLP	4	104044	95.34 \pm 0.138	97.43 \pm 0.470	23.25	22.74 \pm 1.48hr	104380	56.34 \pm 0.181	65.11 \pm 1.685	182.25	29.48 \pm 0.20hr	
GCN	4	101365	90.70 \pm 0.223	97.19 \pm 0.223	127.50	83.41 \pm 0.99hr	101657	55.71 \pm 0.381	69.52 \pm 1.948	142.50	109.70 \pm 4.39hr	
GraphSage	4	101739	86.18 \pm 0.001	90.98 \pm 0.001	98.00	11.00 \pm 1.13hr	101704	59.19 \pm 0.002	59.19 \pm 0.002	124.50	124.50 \pm 0.00hr	
MoNet	4	100449	90.80 \pm 0.032	96.69 \pm 0.440	146.25	93.19 \pm 0.25hr	104229	54.65 \pm 0.518	65.53 \pm 0.715	141.50	97.14 \pm 0.33hr	
GAT	4	100449	95.55 \pm 0.026	99.99 \pm 0.008	104.75	42.26 \pm 0.25hr	110704	64.22 \pm 0.455	89.11 \pm 0.499	103.75	55.27 \pm 1.62hr	
GatedGCN	4	104217	95.34 \pm 0.143	100.00 \pm 0.000	96.25	126.79 \pm 0.25hr	104357	67.31 \pm 0.311	94.55 \pm 1.018	97.00	154.25 \pm 1.26hr	
GIN	4	105434	96.48 \pm 0.052	100.00 \pm 0.000	128.00	39.22 \pm 1.41hr	105654	55.25 \pm 1.527	79.42 \pm 0.970	141.50	52.12 \pm 0.27hr	
RingGNN	2	505182	91.86 \pm 0.049	92.16 \pm 0.005	16.25	257.59 \pm 0.60hr	105168	19.36 \pm 1.108	40.39 \pm 1.397	112.50	112.50 \pm 0.00hr	
	5	506357	Diverged	Diverged	Diverged	510439	Diverged	Diverged	Diverged	Diverged	259.24 \pm 12.60hr	
3WLGNN	3	108924	95.07 \pm 0.061	95.83 \pm 0.138	27.75	152.32 \pm 0.24hr	108516	59.17 \pm 1.593	63.75 \pm 2.697	28.50	160.29 \pm 12.60hr	
	5	508024	95.00 \pm 0.022	95.46 \pm 0.019	26.25	166.87 \pm 0.24hr	102797	58.04 \pm 0.212	61.57 \pm 3.575	20.00	291.22 \pm 15.58hr	
	8	500816	Diverged	Diverged	Diverged	103184	Diverged	Diverged	Diverged	Diverged	Diverged	
LINK PREDICTION												
Model	L	#Param	Test F1 \pm s.d.	TSP		#Epoch	COLLAB		#Param	Test Hits \pm s.d.	#Epoch	Epoch/Total
				Train F1 \pm s.d.	Epoch		Epoch/Total	#Param	Test Hits \pm s.d.	Epoch	Epoch/Total	
MLP	4	169972	0.544 \pm 0.022	0.544 \pm 0.101	164.25	10.54 \pm 2.31hr	39441	20.35 \pm 2.168	29.80 \pm 3.360	147.50	2.09 \pm 0.09hr	
GCN	4	97235	0.630 \pm 0.001	0.631 \pm 0.001	261.00	152.89 \pm 1.11hr	40479	9.42 \pm 2.131	9.12 \pm 0.698	122.50	122.50 \pm 0.00hr	
GraphSage	4	101739	0.630 \pm 0.001	0.631 \pm 0.001	260.00	151.00 \pm 1.13hr	408692	9.18 \pm 0.323	9.18 \pm 0.323	127.50	127.50 \pm 0.00hr	
MoNet	4	99007	0.641 \pm 0.002	0.643 \pm 0.002	282.00	84.46 \pm 0.65hr	39751	36.14 \pm 2.219	61.15 \pm 3.973	167.50	26.69 \pm 1.36hr	
GAT	4	96182	0.671 \pm 0.002	0.673 \pm 0.002	328.25	68.23 \pm 0.25hr	42637	51.30 \pm 0.962	51.85 \pm 1.114	157.00	15.12 \pm 0.80hr	
GatedGCN	4	97958	0.79 \pm 0.001	0.99 \pm 0.001	159.00	40.65 \pm 0.25hr	40965	52.63 \pm 1.168	96.10 \pm 1.876	95.00	45.47 \pm 0.27hr	
GatedGCN-E	4	97858	0.88 \pm 0.003	0.811 \pm 0.003	197.00	218.51 \pm 0.24hr	41899	52.84 \pm 1.345	96.16 \pm 0.453	94.75	45.27 \pm 0.20hr	
GatedGCN-E-PE	16	100570	0.88 \pm 0.003	0.850 \pm 0.003	197.00	149.00 \pm 0.24hr	40965	49.21 \pm 1.560	88.47 \pm 1.858	147.50	2.09 \pm 0.09hr	
GIN	4	99002	0.656 \pm 0.000	0.660 \pm 0.000	273.50	72.73 \pm 5.56hr	39544	41.73 \pm 2.284	70.55 \pm 4.444	140.25	8.66 \pm 0.3hr	
RingGNN	2	106862	0.643 \pm 0.000	0.644 \pm 0.002	2.00	178.95 \pm 17.19hr	39803	-	0.00	-	-	
	5	506564	Diverged	Diverged	Diverged	307000	128.98 \pm 1.00hr	39803	-	0.00	-	-
3WLGNN	3	103566	0.694 \pm 0.073	0.695 \pm 0.073	2.00	17468.1 \pm 16.59hr	39803	-	0.00	-	-	
	5	508681	0.726 \pm 0.311	0.290 \pm 0.312	2.00	17190.7 \pm 16.51hr	39803	-	0.00	-	-	
k-NN Heuristic Matrix Face	0	$k=2$	Test F1: 0.693	-	-	-	-	6054651	44.20 \pm 0.452	100.00 \pm 0.000	254.33	2.66 \pm 0.21hr
GRAPH REGRESSION - ZINC												
Model	L	#Param	Test MAE \pm s.d.	MAE		#Epoch	Epoch/Total		#Param	Test MAE \pm s.d.	#Epoch	Epoch/Total
				Train MAE \pm s.d.	Epoch		Epoch/Total	#Param	Test MAE \pm s.d.	Epoch	Epoch/Total	
MLP	4	108975	0.706 \pm 0.006	0.644 \pm 0.005	116.75	116.00 \pm 0.01hr	104229	-	-	-	-	-
GCN	4	100023	0.706 \pm 0.006	0.644 \pm 0.005	116.75	116.00 \pm 0.01hr	104229	-	-	-	-	-
GraphSage	4	105079	0.367 \pm 0.011	0.128 \pm 0.019	197.00	12.78 \pm 0.71hr	104229	-	-	-	-	-
MoNet	4	106002	0.937 \pm 0.010	0.318 \pm 0.016	188.25	117.90 \pm 0.01hr	104229	-	-	-	-	-
GAT	4	504013	0.920 \pm 0.009	0.093 \pm 0.014	171.75	10.82 \pm 0.52hr	104229	-	-	-	-	-
	16	527383	0.512 \pm 0.003	0.383 \pm 0.024	139.00	32.76 \pm 0.01hr	104229	-	-	-	-	-
GatedGCN	4	105735	0.453 \pm 0.007	0.267 \pm 0.014	173.50	5.76 \pm 0.28hr	104229	-	-	-	-	-
GatedGCN-E	4	105875	0.375 \pm 0.003	0.236 \pm 0.003	194.75	5.37 \pm 0.29hr	104229	-	-	-	-	-
GatedGCN-E-PE	16	105011	0.214 \pm 0.013	0.067 \pm 0.019	185.00	10.70 \pm 0.50hr	104229	-	-	-	-	-
GIN	4	103079	0.387 \pm 0.015	0.139 \pm 0.015	153.25	2.29 \pm 0.10hr	104229	-	-	-	-	-
	5	507459	0.387 \pm									

Benchmarking results

- Result #3 : Anisotropic mechanism improve (isotropic) MP-GNNs.
 - Sparse attention^[1,2] and dense attention^[3] are not injective (unlike f_{GIN}) but experiments showed that they are good at generalization.
 - From a representation view, softmax attention is flexible to represent max/mean/weighted mean w.r.t. contextual information.

- [1] Bahdanau, Cho, Bengio, Neural machine translation by jointly learning to align and translate, 2015
 [2] Velickovic, Cucurull, Casanova, Romero, Lio, Bengio, Graph attention networks, 2017
 [3] Bresson, Laurent, Residual gated graph convnets, 2017

Model	L	#Param	Test Acc. \pm s.d.	PATTERN		#Epoch	NODE CLASSIFICATION		CLUSTER			
				Train Acc. \pm s.d.	Epoch		Epoch/Total	#Param	Test Acc. \pm s.d.	Train Acc. \pm s.d.	#Epoch	
MLP	4	10528	85.48 \pm 0.00	85.19 \pm 0.00	42	8.95 \pm 0.00	106455	53.44 \pm 0.05	20.93 \pm 0.002	42.25	85.72 \pm 0.00	
GCN	4	100023	85.39 \pm 0.04	85.12 \pm 0.13	100	11.80 \pm 0.51hr	501657	53.44 \pm 0.05	34.01 \pm 0.00	20.93 \pm 0.002	76.72 \pm 0.00	
GraphSage	16	100023	81.89 \pm 0.34	78.49 \pm 1.59	81.50	492.19 \pm 1.13hr	501687	68.49 \pm 0.076	71.72 \pm 2.212	79.75	27.028 \pm 0.08hr	
GatedGCN	4	101739	59.16 \pm 0.00	50.47 \pm 0.014	43.75	93.41 \pm 1.17hr	102187	50.45 \pm 0.145	54.37 \pm 0.203	64.00	53.56 \pm 0.97hr	
GatedGCN-E	16	502842	50.49 \pm 0.001	50.48 \pm 0.001	46.50	391.16 \pm 1.19hr	303350	63.84 \pm 0.110	86.71 \pm 0.167	57.75	225.61 \pm 3.70hr	
MoNet	4	103775	85.48 \pm 0.037	85.20 \pm 0.044	89.75	35.71 \pm 0.90hr	104227	58.96 \pm 0.131	58.45 \pm 0.183	76.25	24.04 \pm 0.52hr	
GAT	4	109936	78.24 \pm 1.823	77.88 \pm 1.623	96.00	20.92 \pm 0.57hr	110700	57.73 \pm 0.323	58.33 \pm 0.342	67.25	47.82 \pm 0.00	
GatedGCN	16	526990	78.27 \pm 1.018	90.21 \pm 0.476	53.50	50.33 \pm 0.77hr	527874	70.58 \pm 0.047	76.07 \pm 1.362	73.50	35.94 \pm 0.05hr	
GatedGCN	16	502723	85.58 \pm 0.088	86.07 \pm 0.123	65.25	644.71 \pm 1.91hr	502615	73.84 \pm 0.126	87.88 \pm 0.908	60.00	40.07 \pm 0.81hr	
GatedGCN-PE	16	502457	86.59\pm0.085	86.80 \pm 0.133	65.75	647.94 \pm 1.21hr	502453	76.08\pm0.136	88.91 \pm 0.720	57.75	399.66 \pm 6.58hr	
GIN	4	100884	85.59 \pm 0.011	85.85 \pm 0.030	93.00	15.24 \pm 0.40hr	103544	58.38 \pm 0.236	59.48 \pm 0.337	74.75	10.71 \pm 0.33hr	
RingGNN	2	504766	86.24\pm0.025	86.10 \pm 0.021	72.00	95.97 \pm 1.21hr	524202	22.34 \pm 0.000	22.30 \pm 0.000	43.25	501.84 \pm 6.22hr	
3WLGN	3	505749	Diverged	Diverged	Diverged	Diverged	514380	Diverged	Diverged	Diverged	Diverged	
3WLGN	8	502872	85.34 \pm 0.153	84.98 \pm 0.337	16.00	303.29 \pm 0.56hr	502752	57.37 \pm 0.239	57.27 \pm 0.297	51.00	21.52 \pm 0.52hr	
3WLGN	8	508116	85.34 \pm 0.207	85.27 \pm 0.198	81.75	424.23 \pm 0.56hr	508678	55.48 \pm 0.783	55.73 \pm 0.824	66.00	319.98 \pm 5.70hr	
GRAPH CLASSIFICATION												
Model	L	#Param	Test Acc. \pm s.d.	MNIST		#Epoch	CIFAR10		COLLAB			
				Train Acc. \pm s.d.	Epoch		Epoch/Total	#Param	Test Acc. \pm s.d.	Train Acc. \pm s.d.	#Epoch	
MLP	4	104044	95.34 \pm 0.138	97.43 \pm 0.470	232.25	22.74 \pm 1.48hr	104380	56.34 \pm 0.181	65.13 \pm 1.685	185.25	20.67 \pm 0.71hr	
GCN	4	101365	90.70 \pm 0.223	97.19 \pm 0.223	127.50	83.41 \pm 0.99hr	101657	55.71 \pm 0.381	69.52 \pm 1.948	142.50	109.70 \pm 4.39hr	
GraphSage	4	101400	95.35 \pm 0.025	96.00 \pm 0.025	98.00	11.81 \pm 0.13hr	101704	56.34 \pm 0.025	59.19 \pm 0.20hr	124.00	124.00 \pm 0.00hr	
MoNet	4	100449	90.80 \pm 0.032	96.69 \pm 0.440	146.25	93.19 \pm 0.21hr	104229	54.65 \pm 0.518	65.67 \pm 0.215	141.50	97.14 \pm 0.39hr	
GAT	4	110400	95.35 \pm 0.025	99.99 \pm 0.008	104.75	42.26 \pm 0.25hr	110704	64.22 \pm 0.455	89.11 \pm 0.499	103.75	55.27 \pm 1.62hr	
GatedGCN	4	104217	100.00\pm0.000	100.00 \pm 0.000	96.25	126.79 \pm 0.50hr	104357	67.312\pm0.311	94.55 \pm 0.108	97.00	154.15 \pm 4.22hr	
GIN	4	105434	96.48 \pm 0.052	100.00 \pm 0.000	128.00	39.22 \pm 1.41hr	105654	55.25 \pm 1.527	79.41 \pm 0.970	141.50	52.12 \pm 0.27hr	
RingGNN	2	505182	91.86 \pm 0.449	92.16 \pm 0.805	16.25	2575.96 \pm 0.26hr	103168	19.36 \pm 1.108	20.39 \pm 1.397	112.50	299.24 \pm 16.60hr	
3WLGN	8	506357	Diverged	Diverged	Diverged	Diverged	510439	Diverged	Diverged	Diverged	Diverged	
3WLGN	8	308924	95.02 \pm 0.019	95.60 \pm 0.077	26.25	1523.20 \pm 1.24hr	108516	59.17 \pm 1.593	63.75 \pm 2.697	103.75	150.29 \pm 12.60hr	
3WLGN	8	500816	95.02 \pm 0.067	95.60 \pm 0.077	26.25	1608.73 \pm 1.24hr	102779	58.04 \pm 3.512	61.57 \pm 3.575	20.00	299.12 \pm 5.55hr	
LINK PREDICTION												
Model	L	#Param	Test F1 \pm s.d.	TSP		#Epoch	COLLAB		COLLAB			
				Train F1 \pm s.d.	Epoch		Epoch/Total	#Param	Test Hits \pm s.d.	Train Hits \pm s.d.	#Epoch	
MLP	4	10699	0.544 \pm 0.001	0.544 \pm 0.001	164.25	50.15 \pm 2.31hr	39441	20.35 \pm 2.168	29.80 \pm 3.360	147.50	2.09 \pm 0.00hr	
GCN	4	97502	0.544 \pm 0.001	0.631 \pm 0.001	261.00	152.89 \pm 1.15hr	40749	50.42 \pm 1.131	92.12 \pm 0.695	122.50	35.01 \pm 0.24hr	
GraphSage	4	97230	0.544 \pm 0.001	0.630 \pm 0.001	260.00	152.89 \pm 1.15hr	40869	51.88 \pm 0.373	92.12 \pm 0.695	127.50	35.01 \pm 0.24hr	
MoNet	4	99007	0.641 \pm 0.002	0.643 \pm 0.002	282.00	84.46 \pm 0.65hr	39751	36.144 \pm 2.191	61.156 \pm 3.973	167.50	26.69 \pm 1.36hr	
GAT	4	96182	0.671 \pm 0.002	0.673 \pm 0.002	328.25	68.23 \pm 0.25hr	42637	51.301 \pm 0.962	97.851 \pm 1.114	157.00	18.12 \pm 0.80hr	
GatedGCN	4	97858	0.791 \pm 0.003	0.793 \pm 0.003	159.00	18.26 \pm 0.79hr	40965	52.635 \pm 1.168	96.103 \pm 1.876	95.00	453.47 \pm 12.09hr	
GatedGCN	4	97858	0.808 \pm 0.003	0.811 \pm 0.003	197.00	218.51 \pm 0.24hr	41899	52.849\pm1.345	96.165 \pm 0.453	94.75	452.75 \pm 12.08hr	
GatedGCN	4	1005770	0.808 \pm 0.003	0.850 \pm 0.003	53.00	140.15 \pm 0.00hr	40965	49.212 \pm 1.560	88.347 \pm 3.575	141.50	141.50 \pm 0.00hr	
GIN	4	99002	0.656 \pm 0.000	0.660 \pm 0.000	273.50	72.73 \pm 5.56hr	39544	41.730 \pm 2.284	70.55 \pm 4.444	140.25	8.66 \pm 0.38hr	
RingGNN	2	106862	0.644 \pm 0.000	0.644 \pm 0.002	2.00	1785.93 \pm 17.19hr	-	0.00	0.00	0.00	0.00	
3WLGN	8	506564	Diverged	Diverged	Diverged	Diverged	-	0.00	0.00	0.00	0.00	
3WLGN	8	103566	0.695 \pm 0.073	0.695 \pm 0.073	2.00	17468.1 \pm 16.59hr	-	0.00	0.00	0.00	0.00	
3WLGN	8	508681	0.726 \pm 0.311	0.726 \pm 0.312	2.00	17190.7 \pm 16.51hr	-	0.00	0.00	0.00	0.00	
3WLGN	8	508832	OOM	OOM	OOM	OOM	-	0.00	0.00	0.00	0.00	
k-NN Matrix Factor	0	$k=2$	Test F1: 0.693	-	-	-	-	604561	44.206 \pm 0.452	100.000 \pm 0.000	254.33	2.66 \pm 0.21hr
GRAPH REGRESSION - ZINC												
Model	L	#Param	Test MAE \pm s.d.	ZINC		#Epoch	Epoch/Total		Evaluation Metrics:			
				Train MAE \pm s.d.	Epoch		Epoch/Total	#Param	higher is better, except for ZINC			
MLP	4	108975	0.706 \pm 0.006	0.644 \pm 0.005	116.75	116.00 \pm 0.01hr	-	-	• CLUSTER, PATTERN use weighted accuracy w.r.t. the class sizes.			
GCN	4	100023	0.706 \pm 0.006	0.644 \pm 0.005	116.75	116.00 \pm 0.01hr	-	-	• MNIST, CIFAR10 use multi-label classification accuracy.			
GraphSage	16	105079	0.367 \pm 0.011	0.128 \pm 0.019	197.00	12.78 \pm 0.71hr	-	-	• RingGNN and 3WLGN rely on dense tensors which leads to OOM on both GPU and CPU memory.			
MoNet	4	106002	0.937 \pm 0.010	0.318 \pm 0.016	188.25	117.90 \pm 0.01hr	-	-	• GatedGCN uses H100 via the evaluator provided by OGB [?].			
GAT	4	504013	0.922 \pm 0.009	0.093 \pm 0.014	171.75	10.82 \pm 0.52hr	-	-	• GatedGCN uses H100 via the evaluator provided by OGB [?].			
GatedGCN	4	513185	0.344 \pm 0.007	0.067 \pm 0.004	144.00	12.98 \pm 0.53hr	-	-	• GatedGCN uses H100 via the evaluator provided by OGB [?].			
GatedGCN	4	105735	0.453 \pm 0.011	0.287 \pm 0.014	173.50	5.76 \pm 0.28hr	-	-	• GatedGCN uses H100 via the evaluator provided by OGB [?].			
GatedGCN	4	105875	0.375 \pm 0.005	0.236 \pm 0.004	194.75	5.37 \pm 0.29hr	-	-	• GatedGCN uses H100 via the evaluator provided by OGB [?].			
GatedGCN	4	10589	0.369 \pm 0.005	0.236 \pm 0.005	166.25	20.74 \pm 0.01hr	-	-	• GatedGCN uses H100 via the evaluator provided by OGB [?].			
GatedGCN-E	16	505011	0.214\pm0.013	0.067 \pm 0.019	185.00	10.70 \pm 0.50hr	-	-	• GatedGCN-E uses H100 via the evaluator provided by OGB [?].			
GIN	4											

Benchmarking results

- Result #4 : Explicit edge representation enhance link prediction task. We study three variants

- No edge feature/no edge representation :
Isotropic models (like vanilla GCNs^[1])

$$h_i^{\ell+1} = \sigma \left(\sum_{j \in \mathcal{N}_i} W^\ell h_j^\ell \right)$$

- Edge feature/no edge representation :
Anisotropic models (like GAT^[2])

$$h_i^{\ell+1} = \sigma \left(\sum_{j \in \mathcal{N}_i} f_{V^\ell}(h_i^\ell, h_j^\ell) \cdot W^\ell h_j^\ell \right)$$

- Edge feature and edge representation :
Anisotropic models (like GatedGCNs^[3])

$$h_i^{\ell+1} = \sigma \left(\sum_{j \in \mathcal{N}_i} e_{ij}^\ell \cdot W^\ell h_j^\ell \right), \quad e_{ij}^{\ell+1} = f_{V^\ell}(h_i^\ell, h_j^\ell, e_{ij}^\ell),$$

[1] Kipf, Welling, Semi-supervised classification with graph convolutional networks, 2017

[2] Velickovic, Cucurull, Casanova, Romero, Lio, Bengio, Graph attention networks, 2017

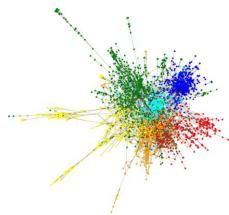
[3] Bresson, Laurent, Residual gated graph convnets, 2017

	Model	E.Feat.	E.Repr.	L	#Param	Test Acc. \pm s.d.	Train Acc. \pm s.d.	#Epochs	Epoch/Total
TSP	GatedGCN	x	x	4	99026	0.646 \pm 0.003	0.648 \pm 0.002	197.50	150.83s/8.34hr
	GatedGCN-E	✓	x	4	98174	0.757 \pm 0.009	0.760 \pm 0.009	218.25	197.80s/12.06hr
		✓	✓	4	97858	0.791 \pm 0.003	0.793 \pm 0.003	159.00	218.20s/9.72hr
	GAT	x	x	4	95462	0.643 \pm 0.001	0.644 \pm 0.001	132.75	325.22s/12.10hr
COLLAB	GAT	✓	x	4	96182	0.671 \pm 0.002	0.673 \pm 0.002	328.25	68.23s/6.25hr
		✓	✓	4	96762	0.748 \pm 0.022	0.749 \pm 0.022	93.00	462.22s/12.10hr
	GAT-E	✓	✓	4	96762	0.782 \pm 0.006	0.783 \pm 0.006	98.00	438.37s/12.11hr
	GatedGCN	x	x	3	26593	35.989 \pm 1.549	60.586 \pm 4.251	148.00	263.62s/10.90hr
COLLAB	GatedGCN	✓	x	3	26715	50.668 \pm 0.291	96.128 \pm 0.576	172.00	384.39s/18.44hr
		✓	✓	3	27055	51.537 \pm 1.038	96.524 \pm 1.704	188.67	376.67s/19.85hr
	GatedGCN-E	✓	✓	3	27055	47.212 \pm 2.016	85.801 \pm 0.984	156.67	377.04s/16.49hr
	GAT	x	x	3	28201	41.141 \pm 0.701	70.344 \pm 1.837	153.50	371.50s/15.97hr
COLLAB	GAT	✓	x	3	28561	50.662 \pm 0.687	96.085 \pm 0.499	174.50	403.52s/19.69hr
		✓	✓	3	26676	49.674 \pm 0.105	92.665 \pm 0.719	201.00	349.19s/19.59hr
COLLAB	GAT-E	✓	✓	3	26676	44.989 \pm 1.395	82.230 \pm 4.941	120.67	328.29s/11.10hr

Outline

- Emerging field
- Tracking progress
- Benchmarking GNNs
 - Architectures
 - Datasets
 - Code infrastructure
 - Experimental setting
 - Results and insights
- **Graph datasets**
- Conclusion

Small real graphs



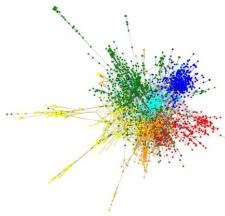
- Small-scale datasets with small-size graphs
 - Cora^[1] : A single graph with 2.7K nodes/5.4K edges, real graph (scientific publications), no feature, node-level classification task (7 classes)
 - Citeseer^[2] : A single graph with 3.3K nodes/4.7K edges, real graph (scientific publications), no feature, node-level classification task (6 classes)
 - TU^[3] : a collection of real (mostly) small-scale datasets
 - TU-ENZYMES : 600 real molecular graphs with 32 nodes/62 edges in average, 18 node features, graph-level classification task (6 classes)
 - TU-PROTEINS : 1.1K real molecular graphs with 39 nodes/72 edges in average, 29 node features, graph-level classification task (2 classes)
 - TU-DD : 1.1K real molecular graphs with 284 nodes/715 edges in average, no feature, graph-level classification task (2 classes)

[1] McCallum, Nigam, Rennie, Seymore, Automating the construction of internet portals with machine learning, 2000
Link to dataset: <https://docs.dgl.ai/generated/dgl.data.CoraGraphDataset.html>

[2] Getoor, Link-based classification, 2005
Link to dataset: <https://docs.dgl.ai/en/0.8.x/generated/dgl.data.CiteseerGraphDataset.html>

[3] Morris et-al, TUDataset: A collection of benchmark datasets for learning with graphs, 2020
Link to datasets: <https://ls11-www.cs.tu-dortmund.de/staff/morris/graphkerneldatasets>

Small real graphs

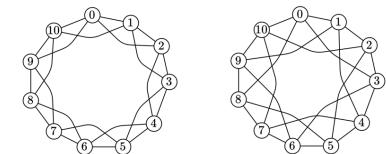


- Please, do not use small datasets as results are not consistent due to high variance^[1] :

	Model	L	#Param	seed 1				seed 2			
				Test Acc. \pm s.d.	Train Acc. \pm s.d.	#Epoch	Epoch/Total	Test Acc. \pm s.d.	Train Acc. \pm s.d.	#Epoch	Epoch/Total
ENZYMES											
ENZYMES	MLP	4	101481	55.833 \pm 3.516	93.062 \pm 7.551	332.30	0.18s/0.17hr	53.833 \pm 4.717	87.854 \pm 10.765	327.80	0.19s/0.18hr
	vanilla GCN	4	103407	65.833\pm4.610	97.688 \pm 3.064	343.00	0.69s/0.67hr	64.833 \pm 7.089	93.042 \pm 4.982	334.30	0.74s/0.70hr
	GraphSage	4	105595	65.000 \pm 4.944	100.000 \pm 0.000	294.20	1.62s/1.34hr	68.167\pm5.449	100.000 \pm 0.000	287.30	1.76s/1.42hr
	MoNet	4	105307	63.000 \pm 8.090	95.229 \pm 5.864	333.70	0.53s/0.49hr	62.167 \pm 4.833	93.562 \pm 5.897	324.40	0.68s/0.62hr
	GAT	4	101274	68.500\pm5.241	100.000 \pm 0.000	299.30	0.70s/0.59hr	68.500\pm4.622	100.000 \pm 0.000	309.10	0.76s/0.66hr
	GatedGCN	4	103409	65.667\pm4.899	99.979 \pm 0.062	316.80	2.31s/2.05hr	70.000\pm4.944	99.979 \pm 0.062	313.20	2.63s/2.30hr
	GIN	4	104864	65.333 \pm 6.823	100.000 \pm 0.000	402.10	0.53s/0.61hr	67.667 \pm 5.831	100.000 \pm 0.000	404.90	0.60s/0.68hr
DD	RingGNN	2	103538	18.667 \pm 1.795	20.104 \pm 2.166	337.30	7.12s/6.71hr	45.333 \pm 4.522	56.792 \pm 6.081	497.50	8.05s/11.16hr
	3WLGNN	3	104658	61.000 \pm 6.799	98.875 \pm 1.571	381.80	9.22s/9.83hr	57.667 \pm 9.522	96.729 \pm 5.525	336.50	11.80s/11.09hr
	MLP	4	100447	72.239 \pm 3.854	73.816 \pm 1.015	371.80	6.36s/6.61hr	72.408\pm3.449	73.880 \pm 0.623	349.60	1.13s/1.11hr
	vanilla GCN	4	102293	72.758 \pm 4.083	100.000 \pm 0.000	266.70	3.56s/2.66hr	73.168\pm5.000	100.000 \pm 0.000	270.20	3.81s/2.88hr
PROTEINS	GraphSage	4	102577	73.433\pm3.429	100.000 \pm 0.000	267.20	11.50s/8.59hr	71.900 \pm 3.647	100.000 \pm 0.000	265.50	6.60s/4.90hr
	MoNet	4	102305	71.736 \pm 3.365	81.003 \pm 2.593	252.60	3.30s/2.34hr	71.479 \pm 2.167	81.268 \pm 2.295	253.50	2.83s/2.01hr
	GAT	4	100132	75.900\pm3.824	95.851 \pm 2.575	201.30	6.31s/3.56hr	74.198\pm3.076	96.964 \pm 1.544	220.10	2.84s/1.75hr
	GatedGCN	4	104165	72.918\pm2.090	82.796 \pm 2.242	300.70	12.05s/10.13hr	71.983 \pm 3.644	83.243 \pm 3.716	323.60	8.78s/7.93hr
	GIN	4	103046	71.910 \pm 3.873	99.851 \pm 0.136	275.70	5.28s/4.08hr	70.883 \pm 2.702	99.883 \pm 0.088	276.90	2.31s/1.79hr
	RingGNN	2	109857	OOD	OOD	OOD	OOD	OOD	OOD	OOD	OOD
	3WLGNN	3	104124	OOD	OOD	OOD	OOD	OOD	OOD	OOD	OOD
PROTEINS	MLP	4	100643	75.644 \pm 2.681	79.847 \pm 1.551	244.20	0.42s/0.29hr	75.823 \pm 2.915	79.442 \pm 1.443	241.20	0.35s/0.24hr
	vanilla GCN	4	104865	76.098 \pm 2.406	81.387 \pm 2.451	350.90	1.55s/1.53hr	75.912\pm3.064	82.140 \pm 2.706	349.60	1.46s/1.42hr
	GraphSage	4	101928	75.289 \pm 2.419	85.827 \pm 0.839	245.40	3.36s/2.30hr	75.559 \pm 1.907	85.118 \pm 1.171	244.40	3.44s/2.35hr
	MoNet	4	103858	76.452\pm2.898	78.206 \pm 0.548	306.80	1.23s/1.06hr	76.453\pm2.892	78.273 \pm 0.695	289.50	1.26s/1.03hr
	GAT	4	102710	76.277\pm2.410	83.186 \pm 2.000	344.60	1.47s/1.42hr	75.557 \pm 3.443	84.253 \pm 2.348	335.10	1.51s/1.41hr
	GatedGCN	4	104855	76.363\pm2.904	79.431 \pm 0.695	293.80	5.03s/4.13hr	76.721\pm3.106	78.689 \pm 0.692	272.80	4.78s/3.64hr
	GIN	4	103854	74.117 \pm 3.357	75.351 \pm 1.267	420.90	1.02s/1.20hr	71.241 \pm 4.921	71.373 \pm 2.835	362.00	1.04s/1.06hr
PROTEINS	RingGNN	2	109036	67.564 \pm 7.551	67.607 \pm 4.401	150.40	28.61s/12.08hr	56.063 \pm 6.301	59.289 \pm 5.560	222.70	19.08s/11.88hr
	3WLGNN	3	105366	61.712 \pm 4.859	62.427 \pm 4.548	211.40	12.82s/7.58hr	64.682 \pm 5.877	65.034 \pm 5.253	200.40	13.05s/7.32hr

Table 16: Performance on the TU datasets with 10-fold cross validation (higher is better). Two runs of all the experiments using the same hyperparameters but different random seeds are shown separately to note the differences in ranking and variation for reproducibility. The top 3 performance scores are highlighted as **First**, **Second**, **Third**.

Artificial graphs



- Small-scale datasets with small-size graphs for isomorphism
 - CSL^[1] : 150 artificial graphs with 41 nodes/164 edges in average, no feature, graph-level classification task (10 classes)
 - EXP^[2] : 600 artificial graphs with 44 nodes/110 edges in average, no feature, graph-level classification task (2 classes)
 - SR25^[3] : 15 artificial graphs with 25 nodes/300 edges in average, no feature, graph-level classification task (15 classes)
 - Medium-scale datasets with medium-size SBM* graphs for pattern matching and clustering
 - PATTERN^[4] : 14k artificial graphs with 117 nodes/4749 edges in average, 1 node feature, node-level classification task (2 classes)
 - CLUSTER^[4] : 12k artificial graphs with 117 nodes/4301 edges in average, 1 node feature, node-level classification task (6 classes)

[1] Murphy et-al, Relational pooling for graph representations, 2019

Link to dataset: https://pytorch-geometric.readthedocs.io/en/latest/modules/torch_geometric/datasets/gnn_benchmark_dataset.html

[2] Abboud et-al, The surprising power of graph neural networks with random node initialization, 2020

Link to dataset: https://github.com/XiaoxinHe/Graph-MLPMixer/blob/b24ae88cd5af68b06a0ccdc8866f3fb4f432550/core/data_utils/exp.py

[3] Balcilar et-al, Breaking the limits of message passing graph neural networks, 2021

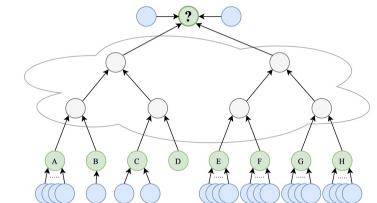
Link to dataset: https://github.com/XiaoxinHe/Graph-MLPMixer/blob/b24ae88cd5af68b06a0ccdc8866f3fb4f432550/core/data_utils/sr25.py

[4] Dwivedi, Joshi, Laurent, Bengio, Bresson, Benchmarking graph neural networks, 2020

Link to datasets: https://pytorch-geometric.readthedocs.io/en/latest/modules/torch_geometric/datasets/gnn_benchmark_dataset.html

* SBM = Stochastic Block Models

Artificial graphs



- Medium-scale datasets with medium-size graphs for graph theory
 - Cycles^[1] : 20k artificial graphs with 48 nodes/87 edges in average, no feature, graph-level classification task (2 classes)
 - GraphTheoryProp^[2] : 7k artificial graphs with 18 nodes/95 edges in average, no feature, collection of 6 graph-level and node-level tasks
 - TSP^[3] : 12k artificial random graphs with 275 nodes/6894 edges in average, 3 node features, 1 edge feature, link-level classification task (2 classes)
- Medium-scale datasets with medium-size graphs for over-squashing
 - TreeNeighbourMatch^[4] : up to 640k artificial tree-graphs with 511 nodes/510 edges in average, no feature, node-level classification task (256 classes)
- Large-scale datasets with small-size graphs
 - GraphWorld^[5] : up to 1M SBM artificial graphs, up to 512 nodes, collection of graph-level, link-level node-level tasks

[1] Loukas, What graph neural networks cannot learn: depth vs width, 2020

Link to dataset: <https://github.com/graphdeeplearning/benchmarking-gnns/blob/master/data/cycles.py>

[2] Corso et-al, Principal neighbourhood aggregation for graph nets, 2020

Link to dataset: <https://github.com/graphdeeplearning/benchmarking-gnns/blob/master/data/graphtheoryprop.py>

[3] Dwivedi, Joshi, Laurent, Bengio, Bresson, Benchmarking graph neural networks, 2020

Link to dataset: https://pytorch-geometric.readthedocs.io/en/latest/_modules/torch_geometric/datasets/gnn_benchmark_dataset.html

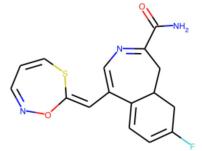
[4] Alon, Yahav, On the bottleneck of graph neural networks and its practical implications, 2020

Link to dataset: https://github.com/XiaoxinHe/Graph-MLPMixer/blob/b24ae88cd5af68b06a0ccdc8866f3fbb4f432550/core/data_utils/tree_dataset.py

[5] Palowitch et-al, Graphworld: Fake graphs bring real insights for gnns, 2022

Link to dataset: <https://github.com/google-research/graphworld>

Molecular graphs



- Medium-scale datasets with small-size graphs
 - QM9^[1] : 130k real molecular graphs with 9 nodes/19 edges in average, 9 node features, 2 edge features, collection of 19 graph-level regression tasks
 - ZINC^[2] : 12k (and 250k) real molecular graphs with 23 nodes/24 edges in average, 1 node feature, 1 edge feature, graph-level regression task (synthetic target)
 - AQSOL^[2] : 9.8k real molecular graphs with 17 nodes/35 edges in average, 1 node feature, 1 edge feature, graph-level regression task (real target)
 - OGBG-MolTOX21^[3] : 7.8k real molecular graphs with 18 nodes/38 edges in average, 9 node feature, 3 edge feature, collection of 12 graph-level classification tasks (2 classes)
 - OGBG-MolHIV^[3] : 41k real molecular graphs with 25 nodes/27 edges in average, 9 node feature, 3 edge feature, graph-level classification task (2 classes)
 - OGBG-PPA^[3] : 158k real molecular graphs with 243 nodes/2266 edges in average, 9 node feature, 3 edge feature, graph-level classification task (37 classes)
 - OGBG-MolPCBA^[3] : 437k real molecular graphs with 26 nodes/28 edges in average, 9 node feature, 3 edge feature, collection of 128 graph-level classification tasks (2 classes)

[1] Wu et-al, MoleculeNet: a benchmark for molecular machine learning, 2018

Link to dataset: https://pytorch-geometric.readthedocs.io/en/latest/generated/torch_geometric.datasets.QM9.html#torch_geometric.datasets.QM9

[2] Dwivedi, Joshi, Laurent, Bengio, Bresson, Benchmarking graph neural networks, 2020

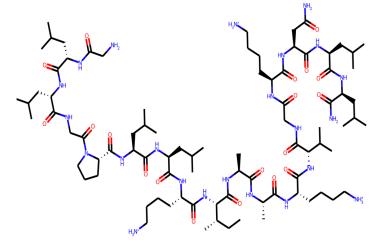
Link to ZINC dataset: https://pytorch-geometric.readthedocs.io/en/latest/generated/torch_geometric.datasets.ZINC.html#torch_geometric.datasets.ZINC

Link to AQSOL dataset: https://github.com/graphdeeplearning/benchmarking-gnns/blob/master/data/molecules/generate_AqSol_raw.ipynb

[3] Hu et-al, Open graph benchmark: Datasets for machine learning on graphs, 2020

Link to datasets and loaders: <https://ogb.stanford.edu/docs/graphprop/#ogbg-mol>

Molecular graphs



- Medium-scale datasets with medium-size graphs
 - Peptides-func^[1] : 15k real molecular graphs with 150 nodes/307 edges in average, 9 node features, 3 edge feature, graph-level classification task (10 classes)
 - Peptides-struct^[1] : 15k real molecular graphs with 150 nodes/307 edges in average, 9 node features, 3 edge feature, collection of 11 graph-level regression tasks
- Large-scale datasets with small-size graphs
 - OGBG-PCQM4Mv2^[2] : 3.3M real molecular graphs with 14 nodes/14 edges in average, 9 node features, 3 edge features, graph-level regression task

[1] Dwivedi et-al, Long range graph benchmark, 2022

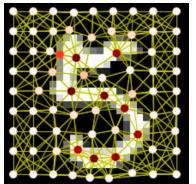
Link to Peptides-func dataset: https://github.com/XiaoxinHe/Graph-MLPMixer/blob/b24ae88cd5af68b06a0ccdc8866f3fb4f432550/core/data_utils/peptides_functional.py

Link to Peptides-struct dataset: https://github.com/XiaoxinHe/Graph-MLPMixer/blob/b24ae88cd5af68b06a0ccdc8866f3fb4f432550/core/data_utils/peptides_structural.py

[2] Wu et-al, MoleculeNet: a benchmark for molecular machine learning, 2018

Link to dataset: <https://ogb.stanford.edu/docs/lsc/pcqm4mv2>

Hybrid real-artificial graphs



- Medium-scale datasets with small-size graphs for computer vision
 - MNIST^[1] : 70k superpixel graphs with 70 nodes/684 edges in average, 3 node features, 1 edge feature, graph-level classification task (10 classes)
 - CIFAR^[1] : 60k superpixel graphs with 117 nodes/1129 edges in average, 5 node features, 1 edge feature, graph-level classification task (10 classes)
- Medium-scale datasets with medium-size graphs for computer vision
 - PascalVOC-SP^[2] : 11k superpixel graphs with 479 nodes/2710 edges in average, 14 node features, 1 edge feature, node-level classification task (20 classes)
 - COCO-SP^[2] : 123k superpixel graphs with 476 nodes/2693 edges in average, 14 node features, 1 edge feature, node-level classification task (80 classes)

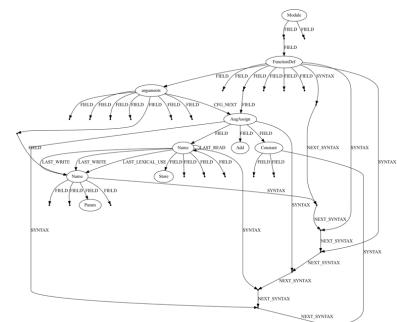
[1] Dwivedi, Joshi, Laurent, Bengio, Bresson, Benchmarking graph neural networks, 2020

Link to datasets: https://pytorch-geometric.readthedocs.io/en/latest/_modules/torch_geometric/datasets/gnn_benchmark_dataset.html

[2] Dwivedi et-al, Long range graph benchmark, 2022

Link to datasets: <https://github.com/vijaydwivedi75/lrgb>

Hybrid real-artificial graphs



- Medium-scale datasets with small-size graphs for codes
 - MalNet-Tiny^[1] : 5k malicious software function call graphs with 1410 nodes/2859 edges in average, no node feature, no edge feature, graph-level classification task (5 classes)
 - OGBG-code2^[2] : 452k abstract syntax graphs of method definitions with 125 nodes/124 edges in average, 4 node features, no edge feature, graph-level classification task (sequence of 5 tokens)
- Large-scale datasets with large-size graphs for codes
 - MalNet^[1] : 1.2M malicious software function call graphs with 15k nodes/35k edges in average, 5 node features, no edge feature, two graph-level classification tasks (47 classes and 696 classes)

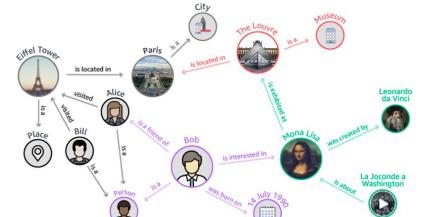
[1] Freitas et-al, A large-scale database for graph representation learning, 2021

Link to dataset: <https://mal-net.org>

[2] Hu et-al, Open graph benchmark: Datasets for machine learning on graphs, 2020

Link to dataset: <https://ogb.stanford.edu/docs/graphprop/#ogbg-code2>

Knowledge graphs



- Single medium-scale graph
 - WikiCS^[1] : a real knowledge graph with 11.7k nodes/216k edges, 300 node feature, no edge feature, node-level classification task (10 classes)
 - FB15k-237^[2] : a real knowledge graph with 14.5k nodes(entities)/310k edges(triplets), 1 node feature (14.5k entity types), 1 edge feature (237 relation types), link-level ranking task
 - WN18^[2] : a real knowledge graph with 40.9k nodes(entities)/142k edges(triplets), 1 node feature (40.9k entity types), 1 edge feature (18 relation types), link-level ranking task
 - OGBL-biokg^[3] : a real knowledge graph with 93k nodes(entities)/5M edges(triplets), 1 node feature (5 entity types), 1 edge feature (51 relation types), link-level ranking task
- Single large-scale graph
 - OGBL-ppa^[3] : a real knowledge graph with 576k nodes(entities)/30M edges(triplets), 1 node feature (58 entity types), no edge feature, link-level ranking task
 - OGBL-wikikg2^[3] : a real knowledge graph with 2.5M nodes(entities)/17M edges(triplets), 1 node feature (2.5M entity types), 1 edge feature (535 relation types), link-level ranking task

[1] Mernyei, Cangea, Wiki-cs: A wikipedia-based benchmark for graph neural networks, 2020

Link to dataset: https://pytorch-geometric.readthedocs.io/en/latest/generated/torch_geometric.datasets.WikiCS.html#torch_geometric.datasets.WikiCS

[2] Bordes et-al, Translating embeddings for modeling multi-relational data, 2013

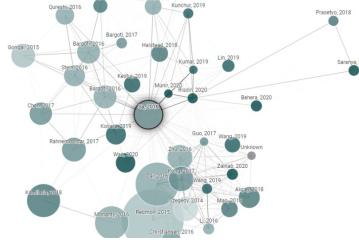
Link to FB15k-237 dataset: https://pytorch-geometric.readthedocs.io/en/latest/generated/torch_geometric.datasets.FB15k_237.html#torch_geometric.datasets.FB15k_237

Link to WN18 dataset: <https://docs.dgl.ai/en/0.9.x/generated/dgl.data.WN18Dataset.html>

[3] Hu et-al, Open graph benchmark: Datasets for machine learning on graphs, 2020

Link to dataset: <https://ogb.stanford.edu/docs/linkprop>

Citation graphs



- Single medium-scale graph
 - OGBN-products^[1] : a real citation graph with 2.5M nodes/61M edges, 100 node features, no edge feature, node-level classification task (47 classes)
 - OGBN-proteins^[1] : a real citation graph with 132k nodes/39M edges, 1 node feature, 8 edge features, collection of 112 node-level classification tasks (2 classes)
 - OGBN-arxiv^[1] : a real citation graph with 169k nodes/1.1M edges, 128 node features, no edge feature, node-level classification task (40 classes)
- Single large-scale graph
 - OGBN-papers100M^[1] : a real citation graph with 111M nodes/1.6G edges, 128 node features, no edge feature, node-level classification task (172 classes)

[1] Hu et-al, Open graph benchmark: Datasets for machine learning on graphs, 2020
Link to dataset: <https://ogb.stanford.edu/docs/nodeprop>

Additional graph datasets

- DGL datasets
 - <https://docs.dgl.ai/api/python/dgl.data.html>
- PyG dataset
 - <https://pytorch-geometric.readthedocs.io/en/latest/modules/datasets.html>
- OGB datasets
 - <https://ogb.stanford.edu>

This screenshot shows the PyG datasets documentation page. It features a sidebar with links for Installation, GET STARTED, TUTORIALS, ADVANCED CONCEPTS, and ADVANCED CONCEPTS. The main content area displays a list of benchmark datasets:

- KarateClub
- TUdataset
- QNNDataset
- Planetoid
- Fakedataset
- FakeParametersDataset
- NELL
- CitationFull
- GeraFull

Each dataset entry includes a brief description and a link to its detailed documentation.

This screenshot shows the DGL datasets documentation page. It has a sidebar with links for GET STARTED, ADVANCED MATERIALS, and API REFERENCE. The main content area is divided into two sections: "Benchmark Datasets" and "Node Prediction Datasets".

Benchmark Datasets

- dgl.data
- CSVdataset

Node Prediction Datasets

- SSTDataset
- KarateClubdataset
- Corafraphdataset
- CiteseerGraphdataset
- PubmedGraphdataset
- CoraFulldataset

This screenshot shows the Open Graph Benchmark homepage. It features a large banner with the text "Open Graph Benchmark" and "Benchmark datasets, data loaders and evaluators for graph machine learning". Below the banner, there are sections for "GET STARTED" and "OGB LSC @ NEURIPS 2022 (NEW)". The page also includes a navigation bar with links for Get Started, Updates, Large-Scale Challenge, Datasets, Leaderboards, Papers, Team, and Github.

The Open Graph Benchmark (OGB) is a collection of realistic, large-scale, and diverse benchmark datasets for machine learning on graphs. OGB datasets are automatically downloaded, processed, and split using the [OGB Data Loader](#). The model performance can be evaluated using the [OGB Evaluator](#) in a unified manner. OGB is a community-driven initiative in active development. We expect the benchmark datasets to evolve. To keep up to date to major updates, subscribe to our google group [here](#).

Lab 1: Loading graph datasets

- Run code01.ipynb

The screenshot shows a Jupyter Notebook interface with the following details:

- Title:** Benchmarking Graph Neural Networks and Graph Datasets
- Section:** Loading graph datasets from DGL, PyG and OGB
- Code Cell (In []):**

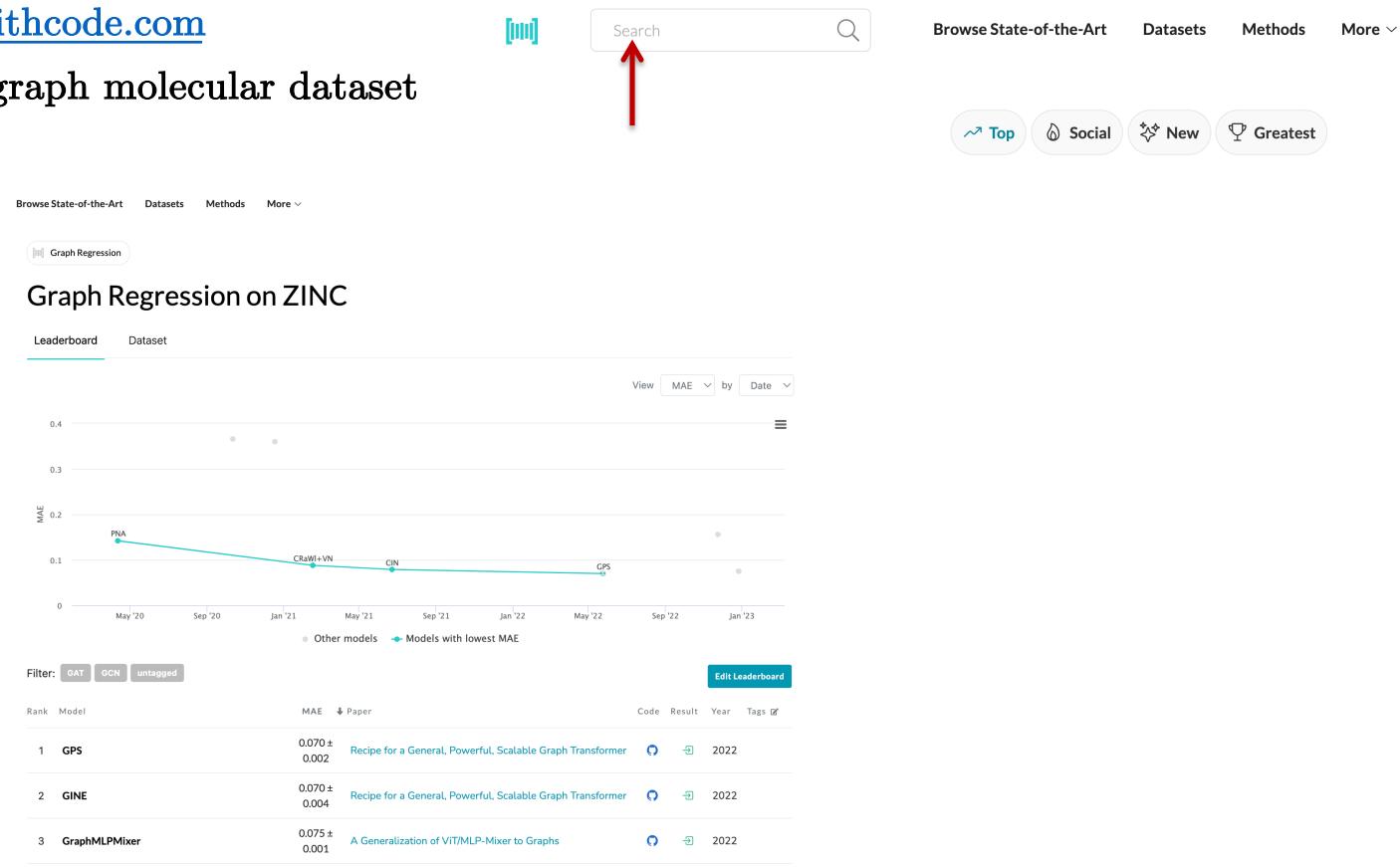
```
# For Google Colaboratory
import sys, os
if 'google.colab' in sys.modules:
    # mount google drive
    from google.colab import drive
    drive.mount('/content/gdrive')
    path_to_file = '/content/gdrive/My Drive/CS6208_codes/codes/labs_lecture05'
    print(path_to_file)
    # change current path to the folder containing "path_to_file"
    os.chdir(path_to_file)
!pwd
!pip install dgl # Install DGL
!pip install torch_geometric # Install PyG
!pip install ogb # Install OGB
```
- Code Cell (In [39]):**

```
import dgl
```
- Code Cell (In [45]):**

```
from dgl.data import CoraGraphDataset
dataset = dgl.data.CoraGraphDataset()
g = dataset[0]
print(g)
num_classes = dataset.num_classes
# get node feature
print(g.ndata['feat'].size())
# get data split
train_mask = g.ndata['train_mask']
val_mask = g.ndata['val_mask']
test_mask = g.ndata['test_mask']
# get labels
print(g.ndata['label'].size())
```

Tracking SOTA per dataset

- Website <https://paperswithcode.com>
 - For example, ZINC graph molecular dataset



Outline

- Emerging field
- Tracking progress
- Benchmarking GNNs
 - Architectures
 - Datasets
 - Code infrastructure
 - Experimental setting
 - Results and insights
- Graph datasets
- Conclusion

Conclusion

- Take-home messages from the benchmark^[1] :
 - MP-GNNs outperformed WL-GNNs on the 8 datasets used in the benchmark.
 - MP-GNNs benefit from graph sparsity, and universal building blocks w/ batch normalization and residual connection.
 - Anisotropic mechanism improves (isotropic) MP-GNNs.
 - Explicit edge representation enhances link prediction.
- Benchmarking bridges the gap between theory and practical performance.
 - However, building good datasets is challenging (size, choice of node/edge features, noise, representativeness, GPU runtime, etc)
 - There is also a limit to learn from a dataset s.a. MNIST and ImageNet. ImageNet has provided many new architectures and insights. But over time, we “overfit” the test accuracy of these datasets with tricks and they do not contribute to new ideas.

[1] Dwivedi, Joshi, Laurent, Bengio, Bresson, Benchmarking graph neural networks, 2020



Questions?