



**ICLR**

---

# **Graph Neural Networks are Inherently Good Generalizers: Insights by Bridging GNNs and MLPs**

**Presenter: Chenxiao Yang**

# Outline

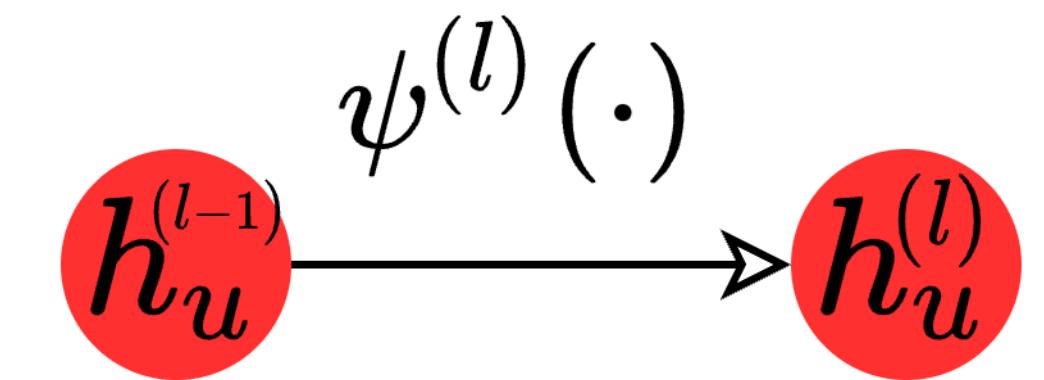
---

- Architectural Connection between GNN and MLP
- Bridging GNN and MLP : Propagational MLP (PMLP)
- Experiments
- Quick Guide of Implementation
- GNNs are Inherently Good Generalizers
- Theoretical Analysis

# Architectural Connection between GNN and MLP

- Each layer of **MLP (multi-layer perceptrons)** consists of linear transformation and non-linear activation :

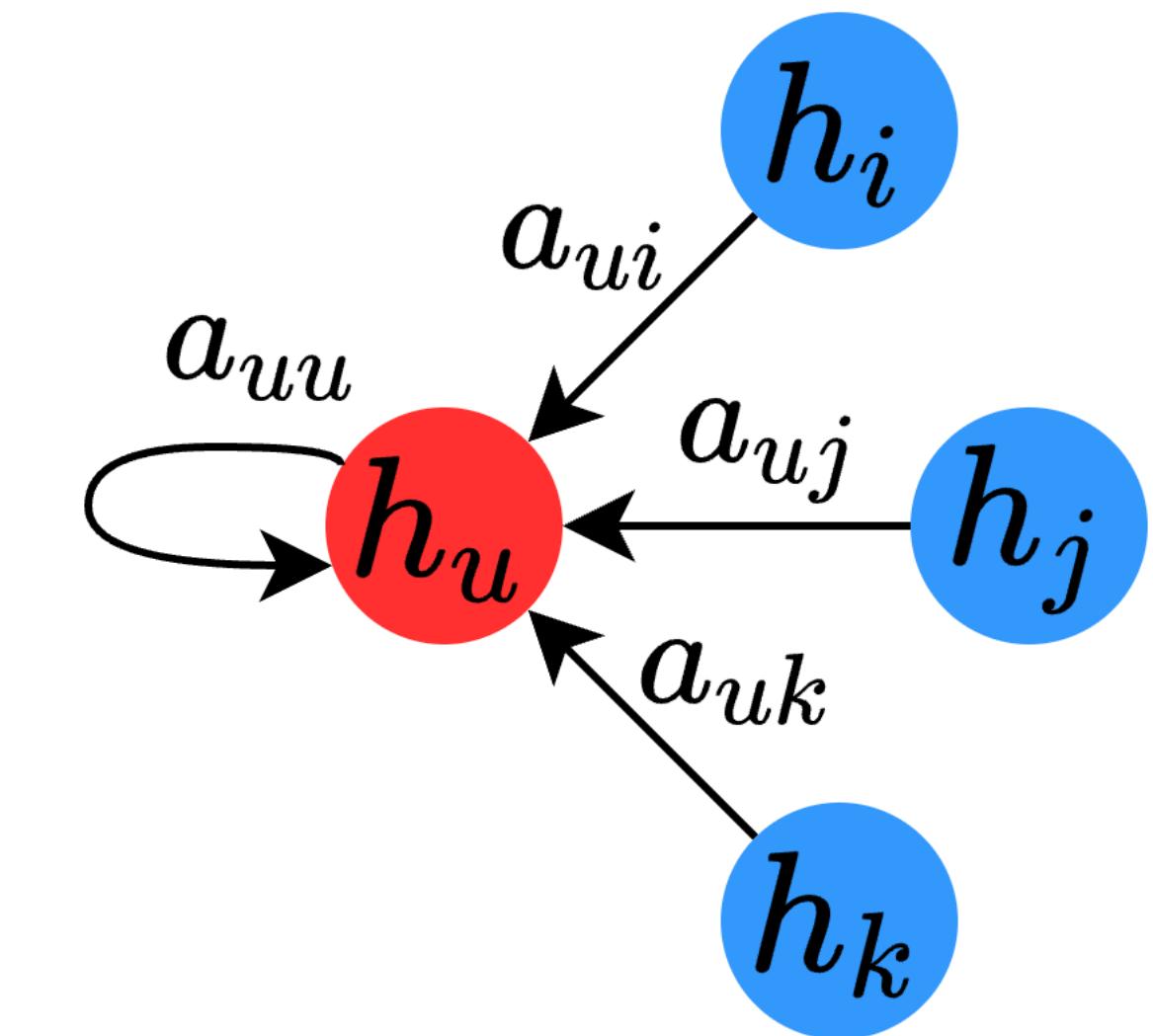
**(Feed-Forward / FF)** :  $\mathbf{h}_u^{(l)} = \psi^{(l)} \left( \mathbf{h}_u^{(l-1)} \right), \quad \psi^{(l)}(\mathbf{h}) = \sigma(\mathbf{h}\mathbf{W}^{(l)})$



- Each layer of **GNN (graph neural networks)** additionally incorporates a message passing operation to propagate node features :

**(Message Passing / MP)** :  $\tilde{\mathbf{h}}_u^{(l-1)} = \sum_{v \in \mathcal{N}_u \cup \{u\}} a_{uv} \cdot \mathbf{h}_v^{(l-1)}$   
e.g.  $a_{uv} = \mathbf{A}_{uv} / \sqrt{\tilde{d}_u \tilde{d}_v}$  for GCN.

**(Feed-Forward / FF)** :  $\mathbf{h}_u^{(l)} = \psi^{(l)} \left( \tilde{\mathbf{h}}_u^{(l-1)} \right), \quad \psi^{(l)}(\mathbf{h}) = \sigma(\mathbf{h}\mathbf{W}^{(l)})$

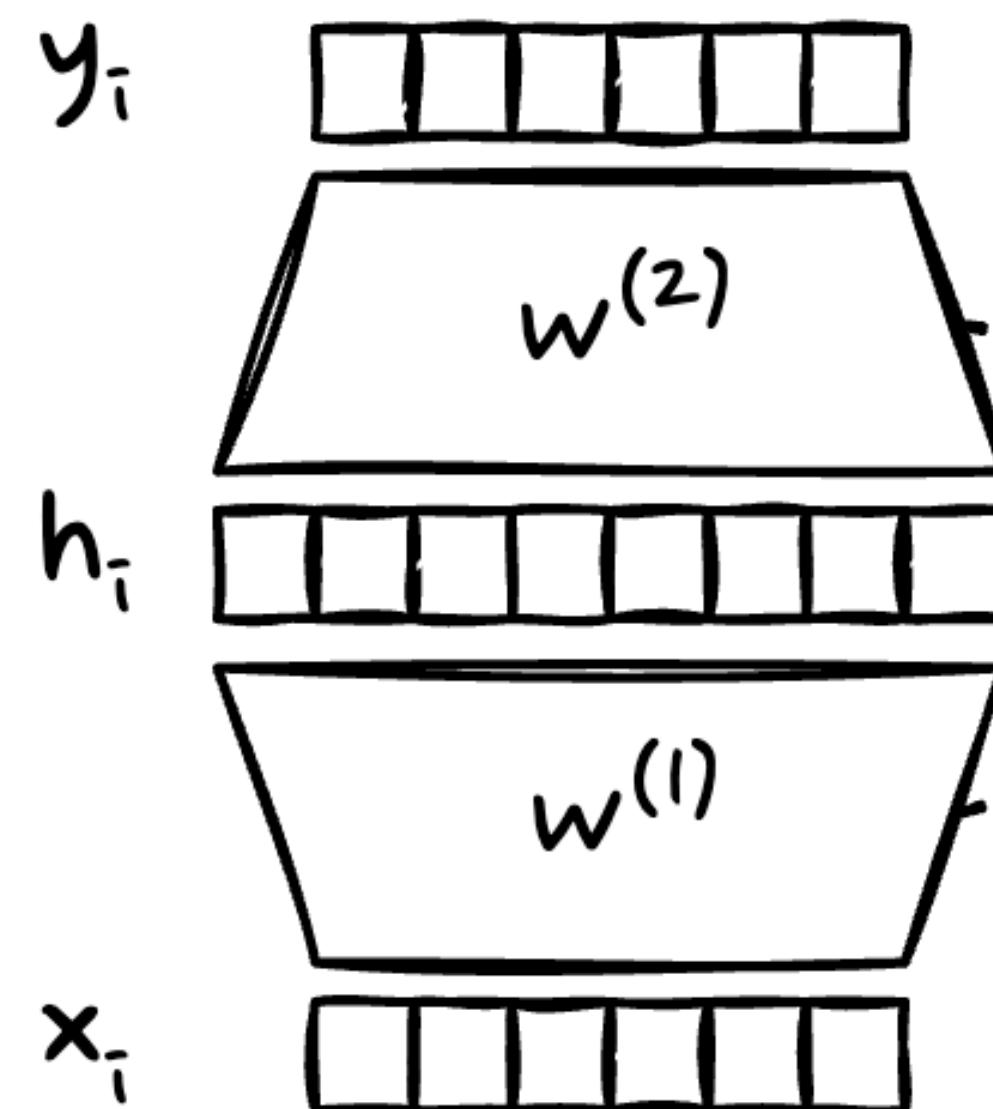


# Architectural Connection between GNN and MLP

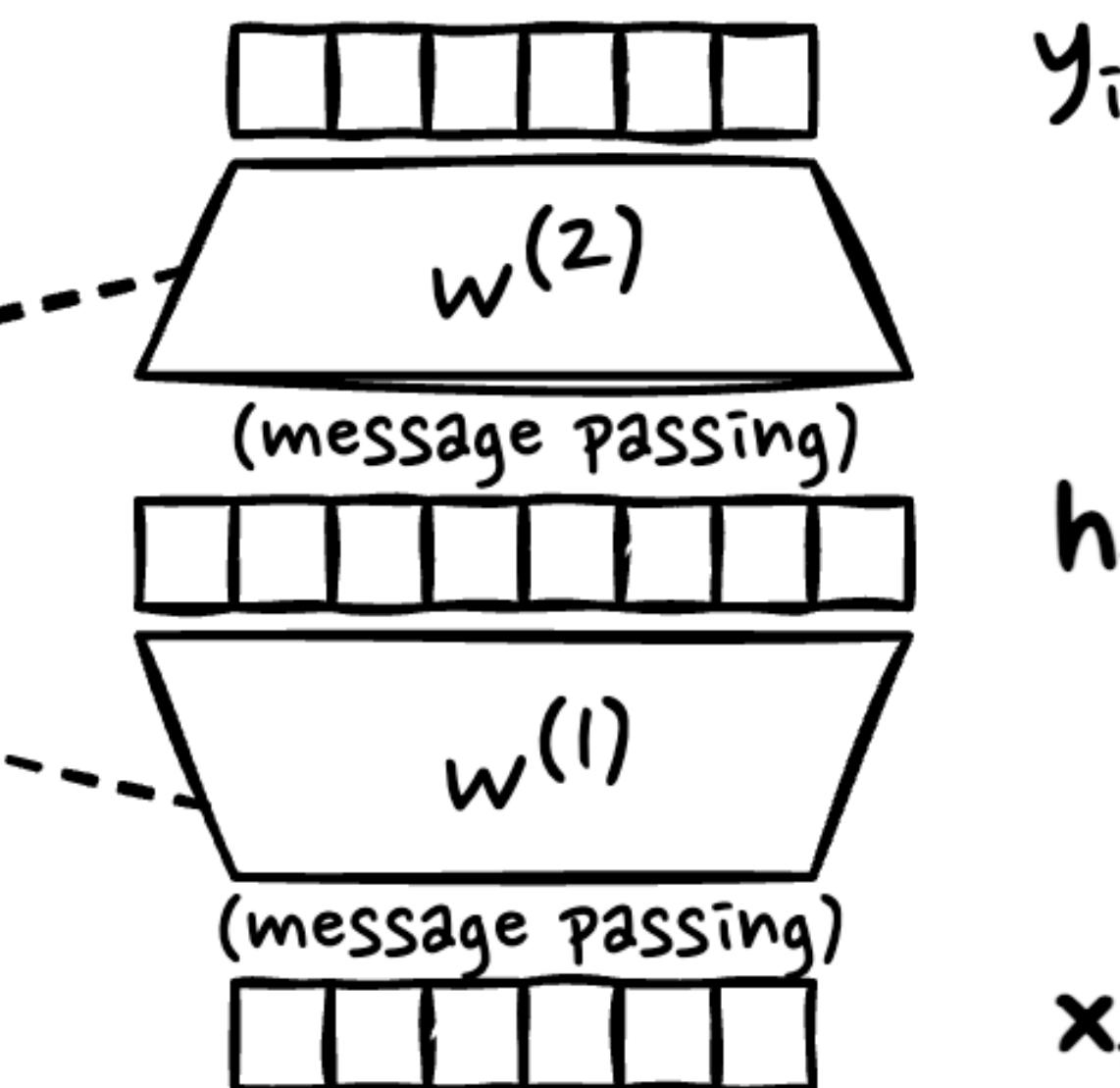
*Q: What if we remove all message passing operations / layers in GNNs?*

- **Observation :** GNNs without message passing operations become MLPs, and they share the same set of learnable parameters.

Multi-Layer Perceptrons



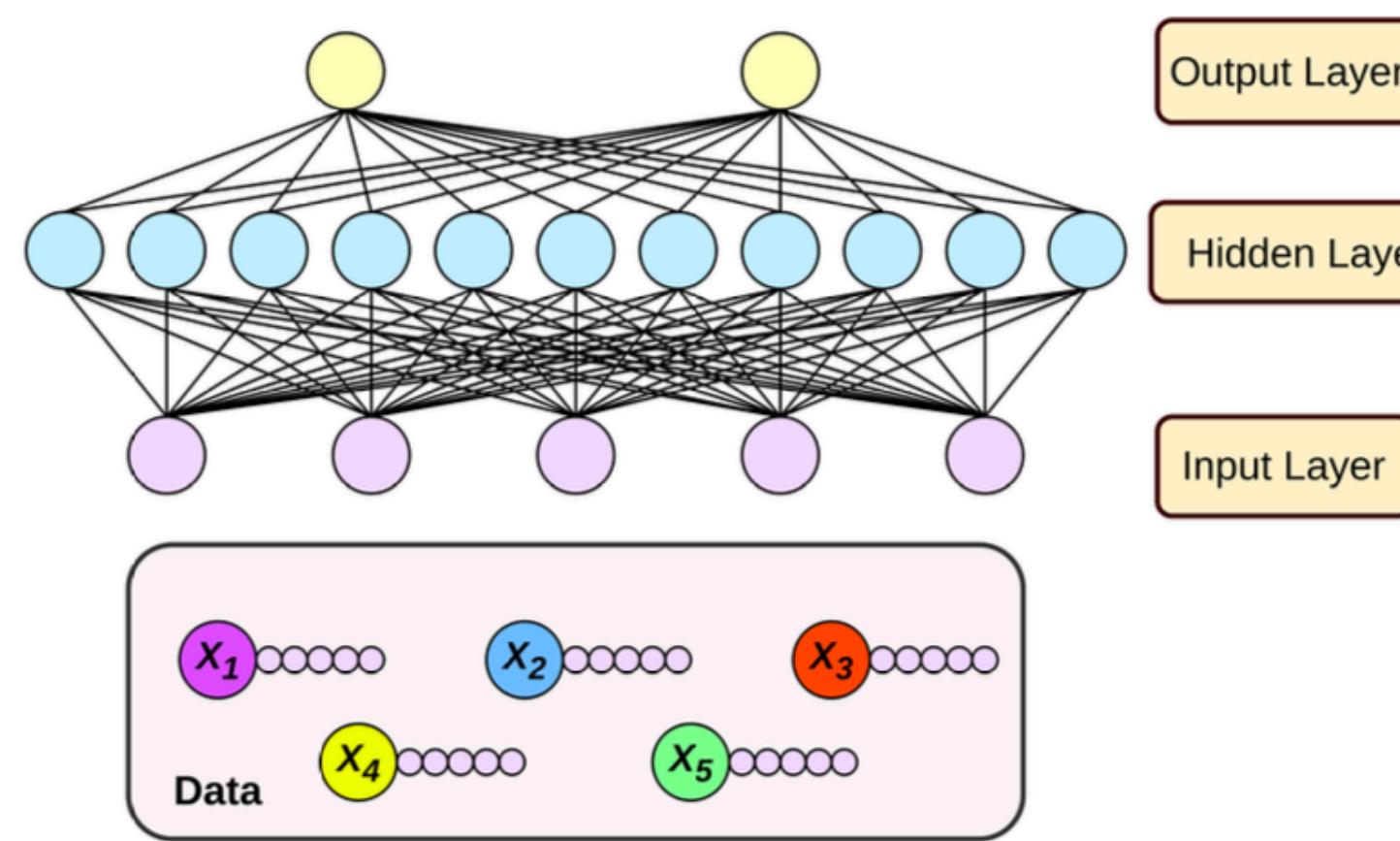
Graph Neural Networks  
(GCN-Style)



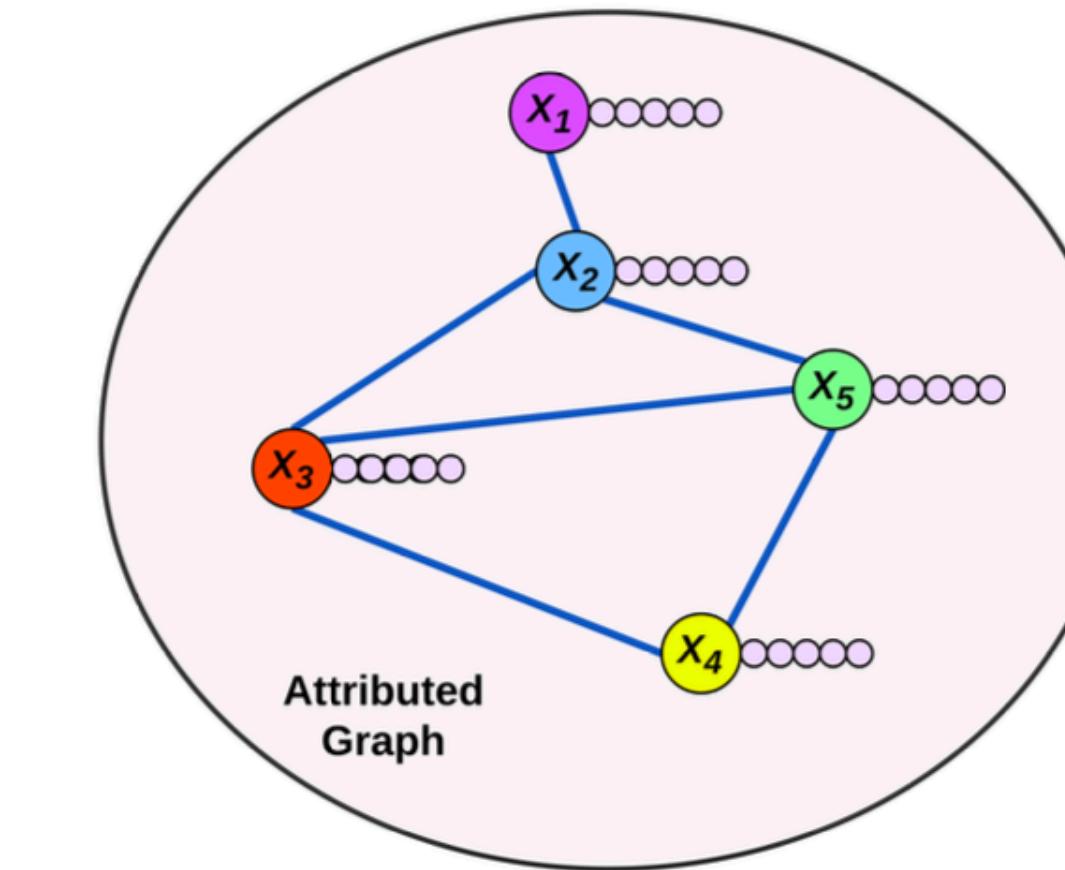
Same Set of  
Parameters

# Architectural Connection between GNN and MLP

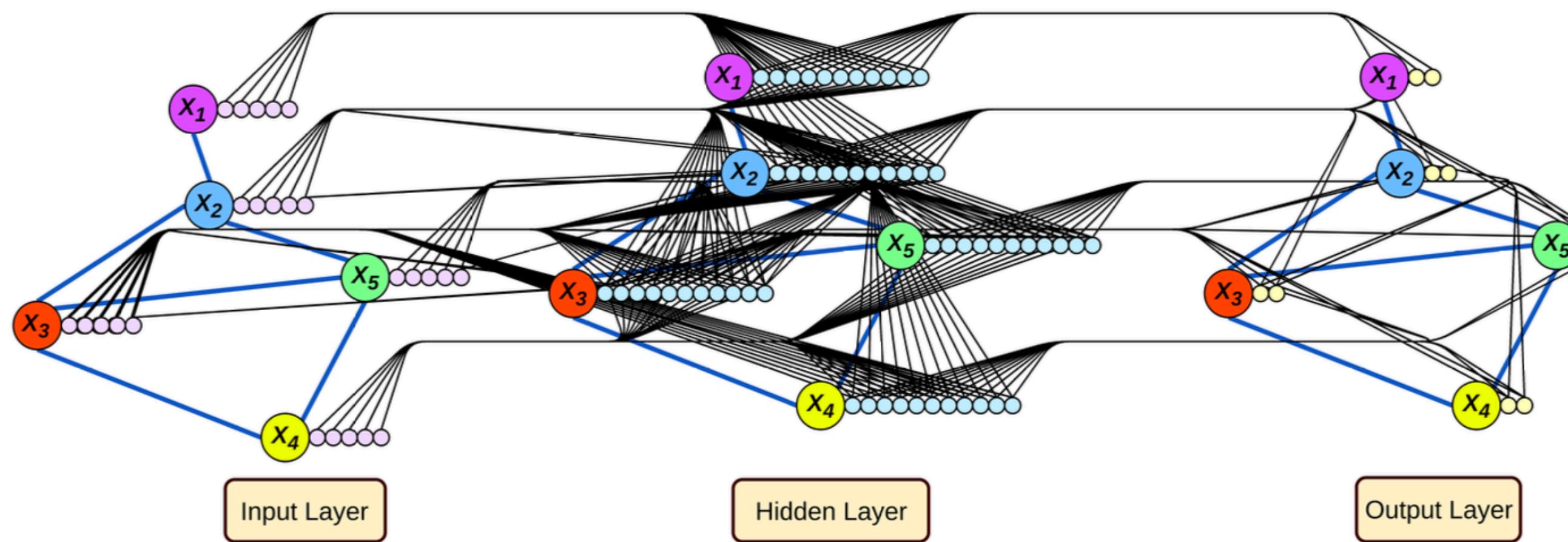
(a) **MLP**



(b) **Graph**



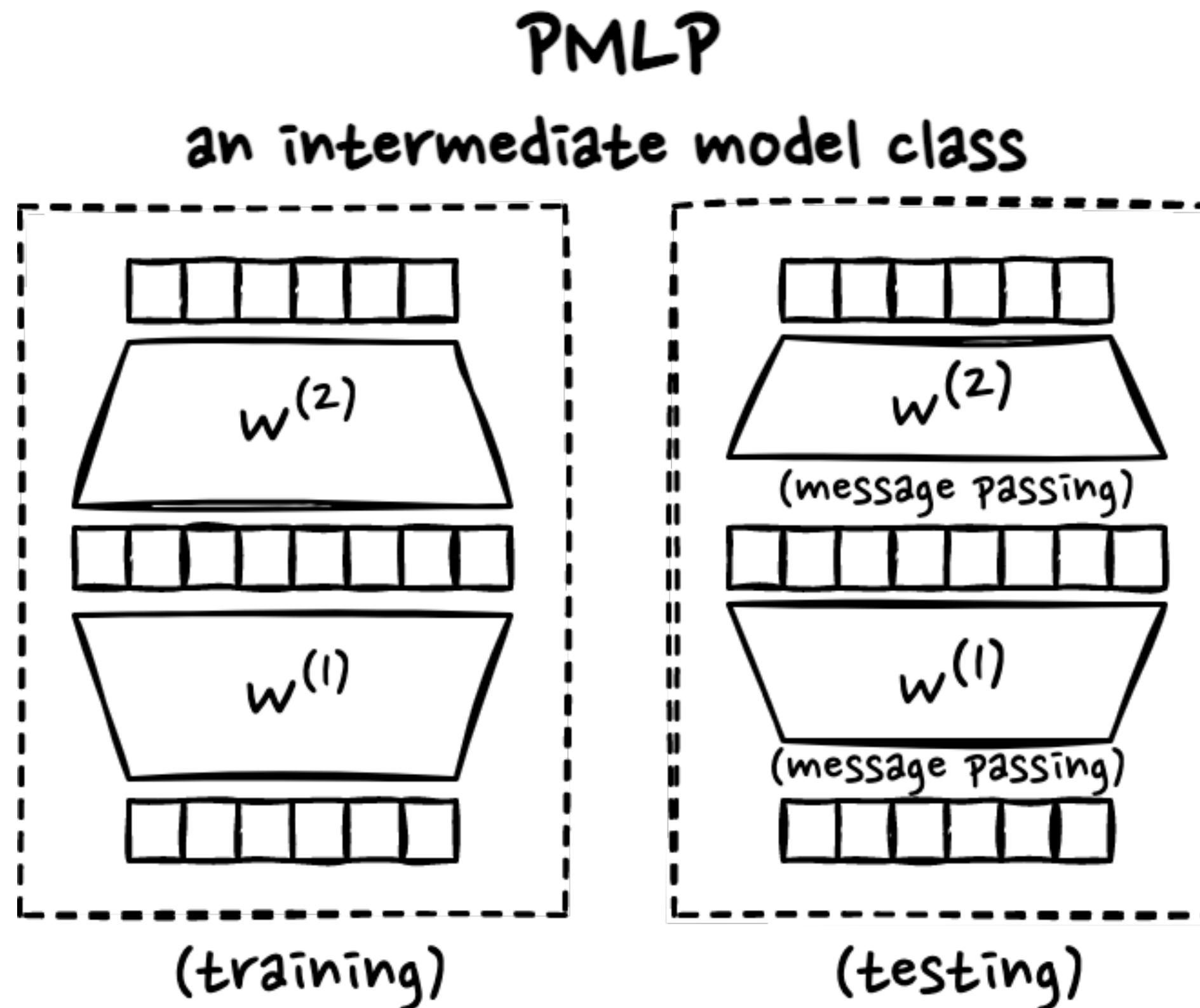
(c) **GNN**



# Bridging GNN and MLP : Propagational MLP (PMLP)

*Q: Does there exist any intermediate models between MLPs and GNNs ?*

$$\text{PMLP} = \text{MLP} \text{ (in training)} + \text{GNN} \text{ (for inference)}$$



- **Training :** PMLPs are exactly the same as a standard MLP (e.g. same architecture, data, loss function, optimization algorithm)
- **Testing / Inference :** PMLPs additionally insert message passing (MP) layers amid feed-forward (FF) layers, to align with various GNN architectures.

# Bridging GNN and MLP : Propagational MLP (PMLP)

---

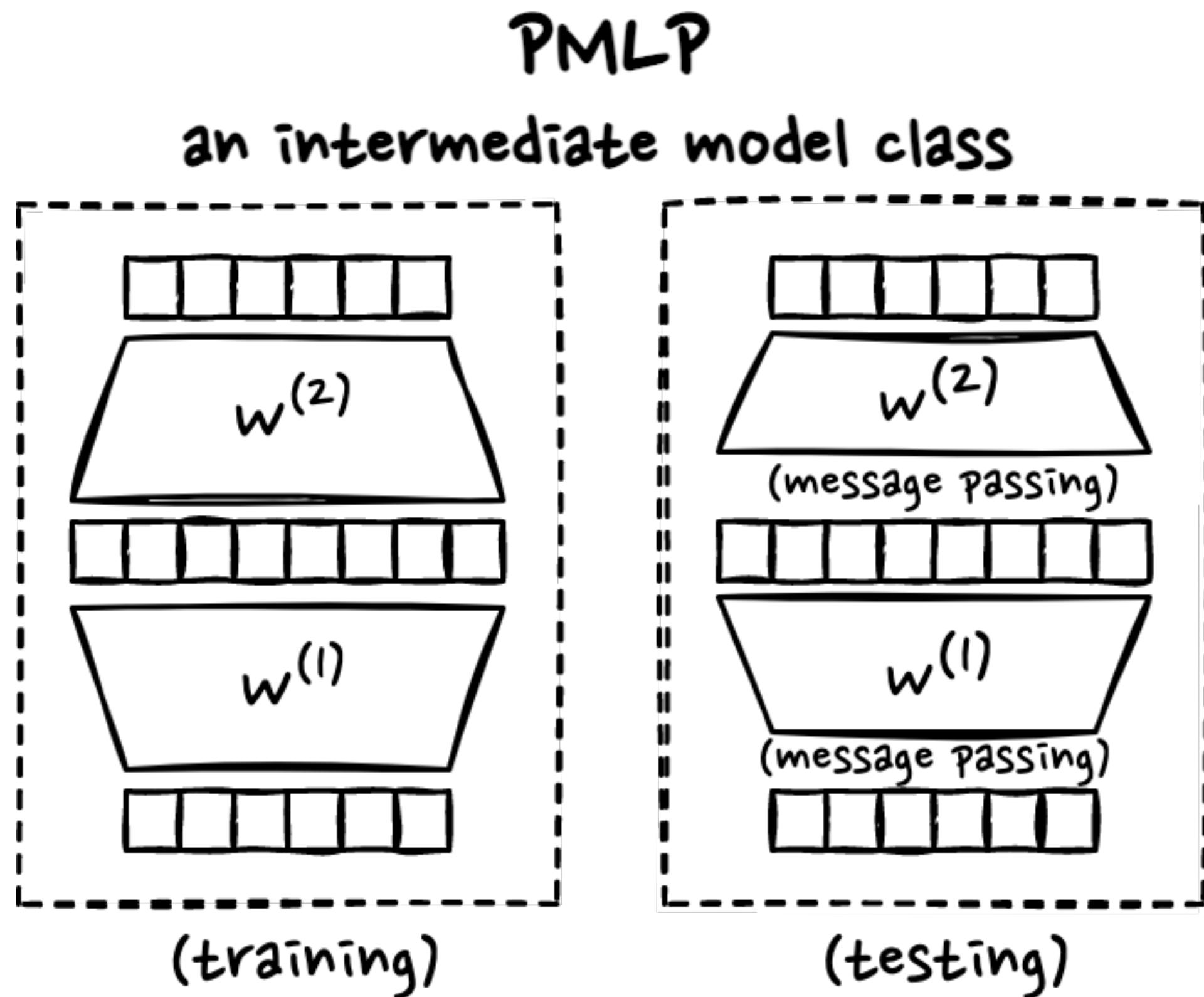
- PMLP has different instantiations depending on the architecture of GNNs.

Model	Train and Valid	Inference
MLP		MLP
PMLP <sub>GCN</sub>	MLP: $\hat{y}_u = \psi(\mathbf{x}_u)$	GCN: $\psi^{(l)}(\text{MP}(\{\mathbf{h}_v^{(l-1)}\}_{v \in \mathcal{N}_u \cup \{u\}}))$
PMLP <sub>SGC</sub>		SGC: $\psi(\text{Multi-MP}(\{\mathbf{x}_v\}_{v \in \mathcal{V}}))$
PMLP <sub>APPNP</sub>		APPNP: Multi-MP( $\psi(\{\mathbf{x}_v\}_{v \in \mathcal{V}})$ )
PMLP <sub>GCNII</sub>	ResNet	GCNII
PMLP <sub>JKNet</sub>	MLP+JK	JKNet

(and more ....)

**Note :** We allow the model after removal of message passing layers to be other models, such as ResNet.

# Bridging GNN and MLP : Propagational MLP (PMLP)



**Equivalent ways to view PMLP :**

1. PMLP = MLP with test-time message passing operations.
2. PMLP = GNN without message passing layers in training.
3. PMLP = Dropping all edges while training GNN.
4. PMLP = Using pretrained weights of MLP on GNN.

# Experiments

*Q: How does PMLP perform compared with MLP and GNN ?*

Dataset		Cora	Citeseer	Pubmed	A-Photo	A-Computer	Coauthor-CS	Coauthor-Physics
#Nodes		2,708	3,327	19,717	7,650	13,752	18,333	34,493
GNNs	GCN	74.82 ± 1.09	67.60 ± 0.96	76.56 ± 0.85	89.69 ± 0.87	78.79 ± 1.62	91.79 ± 0.35	91.22 ± 0.18
	SGC	73.96 ± 0.59	67.34 ± 0.54	76.00 ± 0.59	83.42 ± 2.47	77.10 ± 2.54	91.24 ± 0.59	89.18 ± 0.46
	APPNP	75.02 ± 2.17	66.58 ± 0.77	76.48 ± 0.49	89.51 ± 0.86	78.29 ± 0.55	91.64 ± 0.34	91.80 ± 0.77
MLPs	MLP	55.30 ± 0.58	56.20 ± 1.27	70.76 ± 0.78	75.61 ± 0.63	63.07 ± 1.67	87.51 ± 0.51	85.09 ± 4.11
	$\text{PMLP}_{GCN}$	<b>75.86 ± 0.93</b>	<b>68.00 ± 0.70</b>	<b>76.06 ± 0.55</b>	<b>89.10 ± 0.88</b>	<b>78.05 ± 1.21</b>	<b>91.76 ± 0.27</b>	<b>91.35 ± 0.82</b>
	$\Delta_{GNN}$	+1.39%	+0.59%	-0.65%	-0.66%	-0.94%	-0.03%	+0.14%
	$\Delta_{MLP}$	+37.18%	+21.00%	+7.49%	+17.84%	+23.75%	+4.86%	+7.36%
	$\text{PMLP}_{SGC}$	<b>75.04 ± 0.95</b>	<b>67.66 ± 0.64</b>	<b>76.02 ± 0.57</b>	<b>86.50 ± 1.40</b>	<b>74.72 ± 3.86</b>	<b>91.09 ± 0.50</b>	<b>89.34 ± 1.40</b>
	$\Delta_{GNN}$	+1.46%	+0.48%	+0.03%	+3.69%	-3.09%	-0.16%	+0.18%
	$\Delta_{MLP}$	+35.70%	+20.39%	+7.43%	+14.40%	+18.47%	+4.09%	+4.99%
MLPs	$\text{PMLP}_{APP}$	<b>75.84 ± 1.36</b>	<b>67.52 ± 0.82</b>	<b>76.30 ± 1.44</b>	<b>88.47 ± 1.64</b>	<b>78.07 ± 2.10</b>	<b>91.64 ± 0.46</b>	<b>91.96 ± 0.51</b>
	$\Delta_{GNN}$	+1.09%	+1.41%	-0.24%	-1.16%	-0.28%	+0.00%	+0.17%
	$\Delta_{MLP}$	+37.14%	+20.14%	+7.83%	+17.01%	+23.78%	+4.72%	+8.07%

Table. Experimental results for inductive node classification tasks.

**A vanilla MLP without any special training tricks can be  
as effective as GNNs**

# Experiments

*Q: How does PMLP perform compared with MLP and GNN ?*

	GCN	SGC	APPNP	MLP	<b>PMLP<sub>GCN</sub></b>	<b>PMLP<sub>SGC</sub></b>	<b>PMLP<sub>APP</sub></b>
OGBN-Arxiv ( $\Delta_{GNN}/\Delta_{MLP}$ )	69.04 $\pm$ 0.18	68.56 $\pm$ 0.14	69.19 $\pm$ 0.12	53.86 $\pm$ 0.28	<b>63.74 <math>\pm</math> 2.28</b> (-7.68%/+18.34%)	<b>62.65 <math>\pm</math> 0.35</b> (-8.62%/+16.32%)	<b>63.30 <math>\pm</math> 0.17</b> (-8.51%/+17.53%)
Train loss	1.0480	1.1124	1.0396	1.5917	<b>1.5917</b>	<b>1.5917</b>	<b>1.5917</b>
Train time	63.48 ms	90.50 ms	87.24 ms	7.78 ms	<b>7.78 ms</b>	<b>7.78 ms</b>	<b>7.78 ms</b>
OGBN-Products ( $\Delta_{GNN}/\Delta_{MLP}$ )	71.35 $\pm$ 0.19	71.17 $\pm$ 0.29	70.41 $\pm$ 0.07	56.24 $\pm$ 0.10	<b>69.71 <math>\pm</math> 0.13</b> (-2.30%/+23.95%)	<b>70.09 <math>\pm</math> 0.13</b> (-1.52%/+24.63%)	<b>65.72 <math>\pm</math> 0.09</b> (-6.66%/+16.86%)
Train loss	0.4013	0.4018	0.4311	1.0841	<b>1.0841</b>	<b>1.0841</b>	<b>1.0841</b>
Train time	288.61 ms	665.06 ms	527.26 ms	39.73 ms	<b>39.73 ms</b>	<b>39.73 ms</b>	<b>39.73 ms</b>
Flickr ( $\Delta_{GNN}/\Delta_{MLP}$ )	49.66 $\pm$ 0.57	50.93 $\pm$ 0.16	45.31 $\pm$ 0.28	46.44 $\pm$ 0.14	<b>49.55 <math>\pm</math> 1.30</b> (-0.22%/+6.70%)	<b>50.99 <math>\pm</math> 0.39</b> (+0.12%/+9.80%)	<b>44.31 <math>\pm</math> 0.24</b> (-2.21%/-4.59%)
Train loss	1.1462	1.4030	0.7449	1.3344	<b>1.3344</b>	<b>1.3344</b>	<b>1.3344</b>
Train time	36.82 ms	50.42 ms	163.82 ms	7.44 ms	<b>7.44 ms</b>	<b>7.44 ms</b>	<b>7.44 ms</b>

Table. Experimental results for inductive node classification tasks.

**PMLPs do not need graph information in training, and can reduce the training time up to 65 times**

# Experiments

---

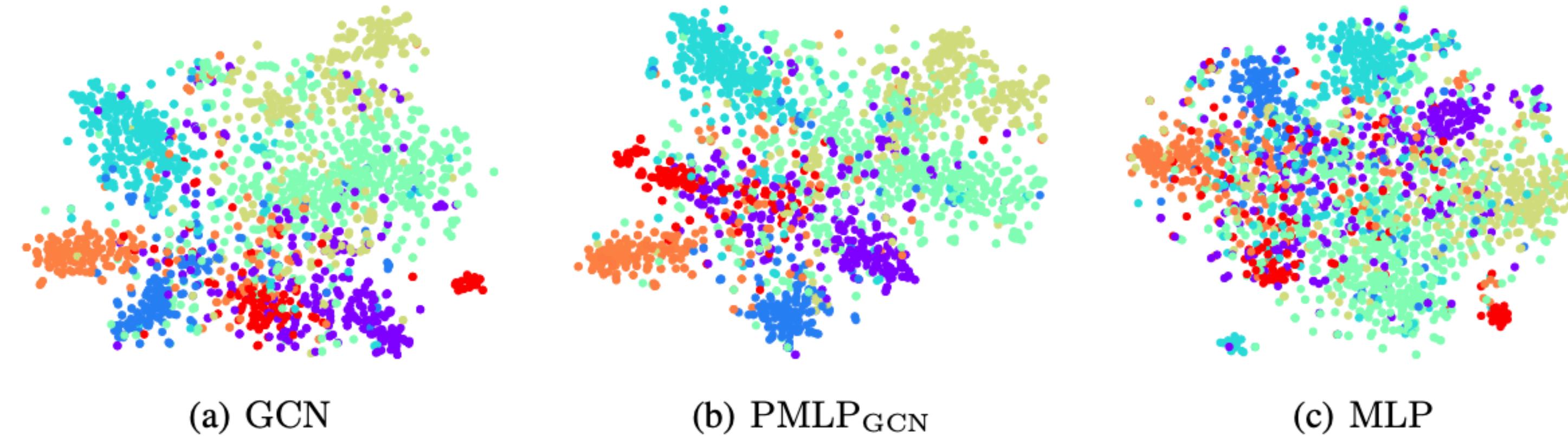


Figure 11: Visualization of node embeddings (2-D projection by t-SNE) in the internal layer for two-layer MLP, GCN and PMLP on Cora.

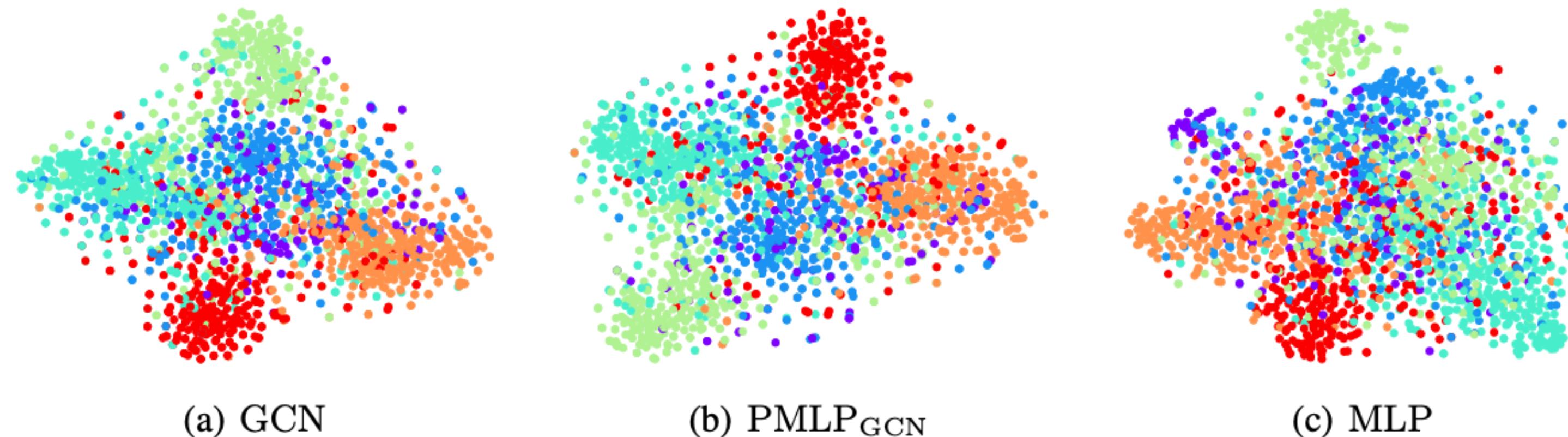
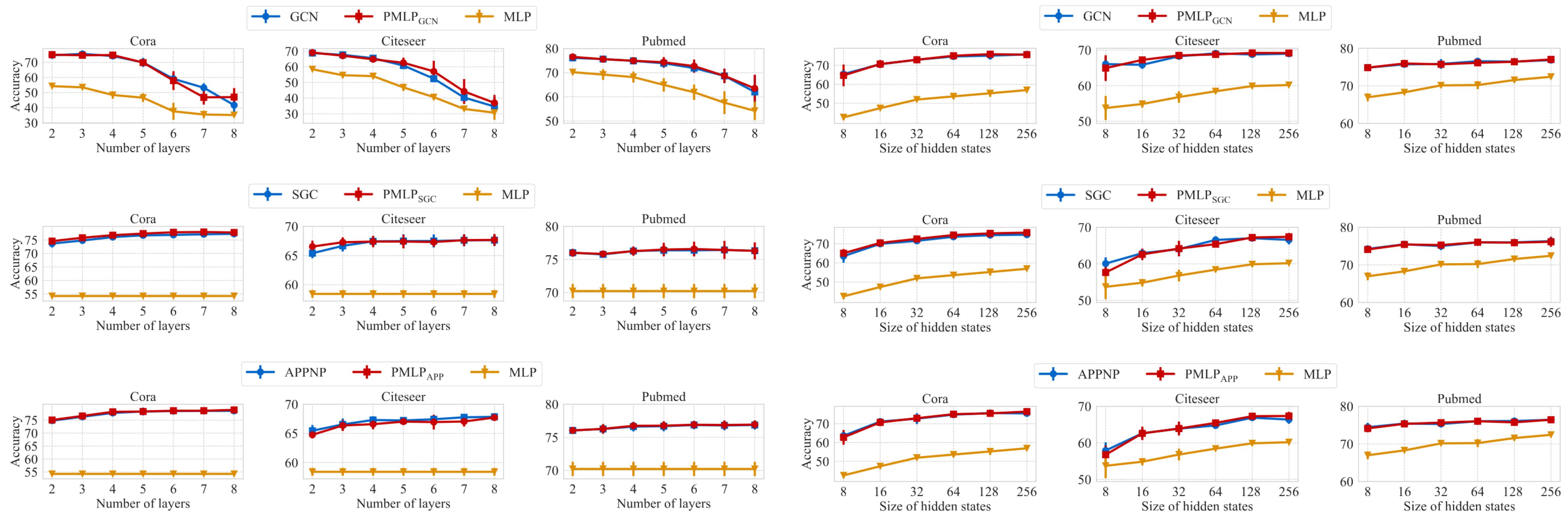


Figure 13: Visualization of node embeddings (2-D projection by t-SNE) in the internal layer for two-layer MLP, GCN and PMLP on Citeseer.

# Experiments

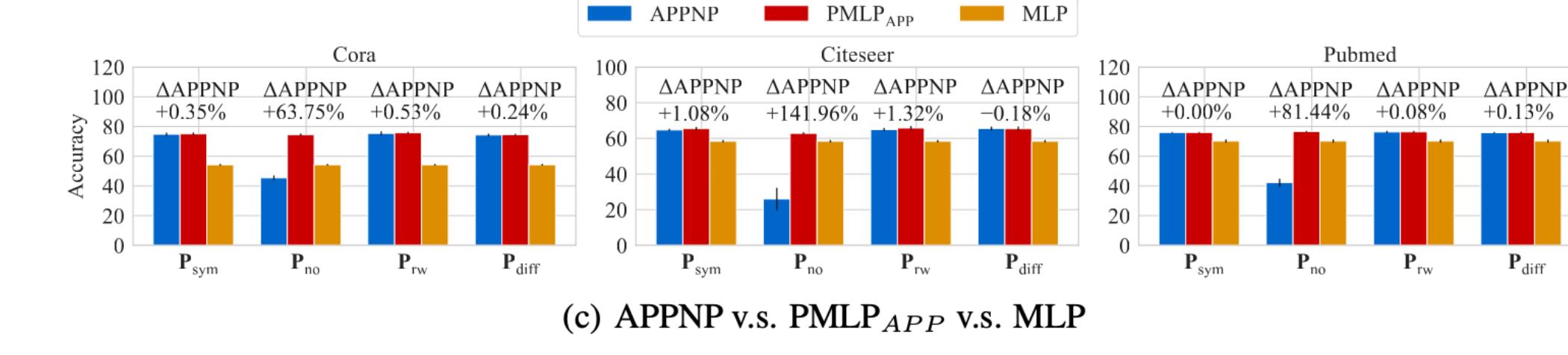
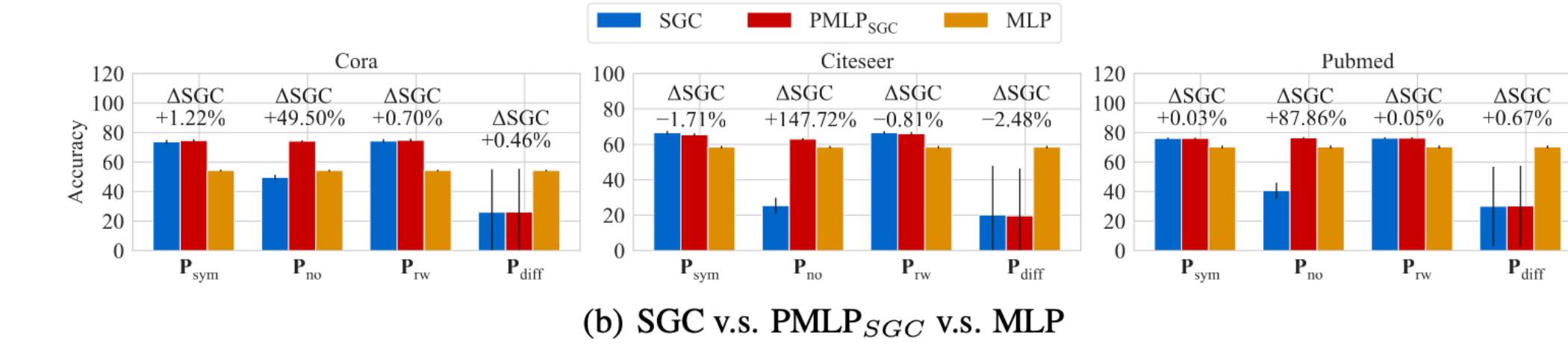
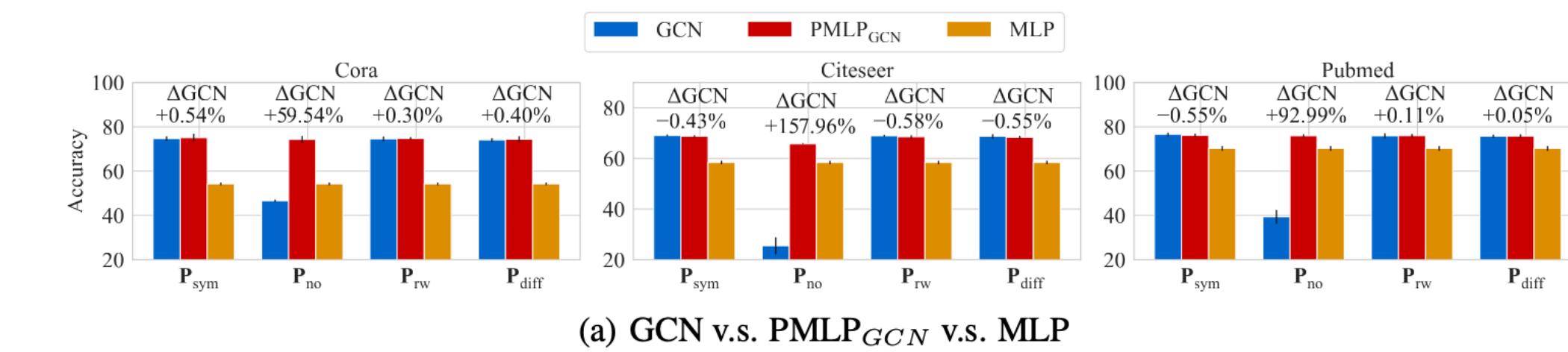
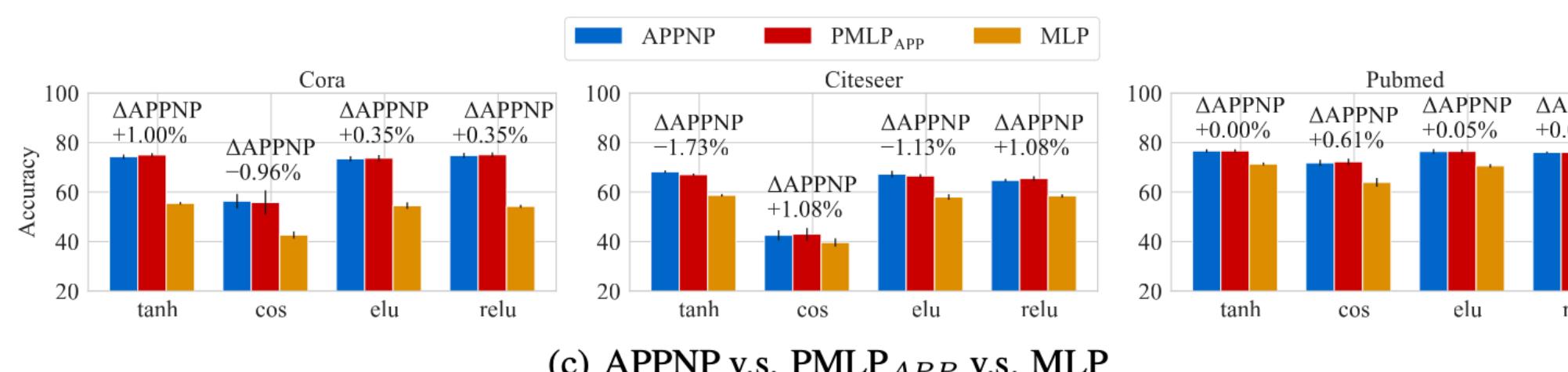
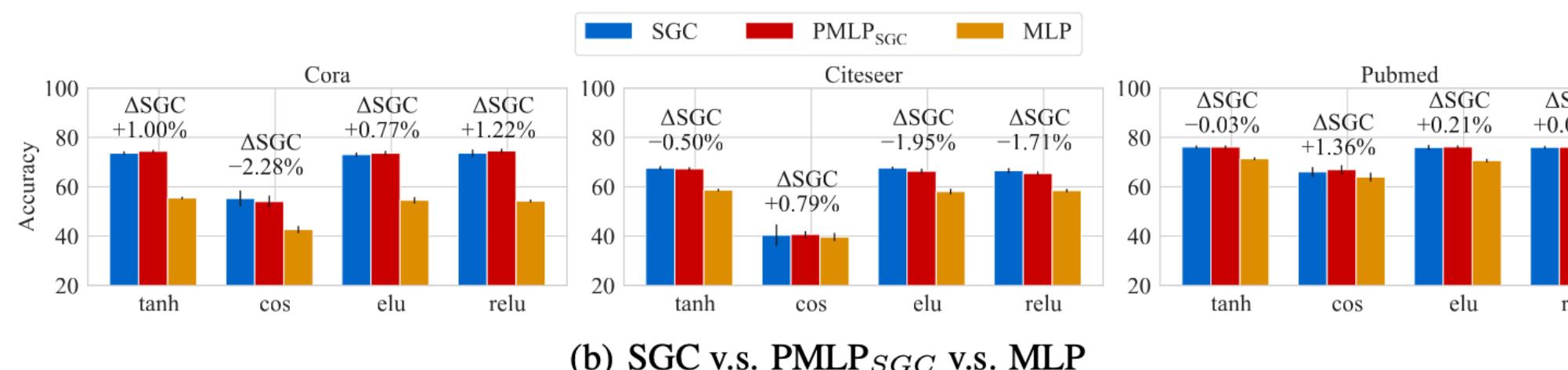
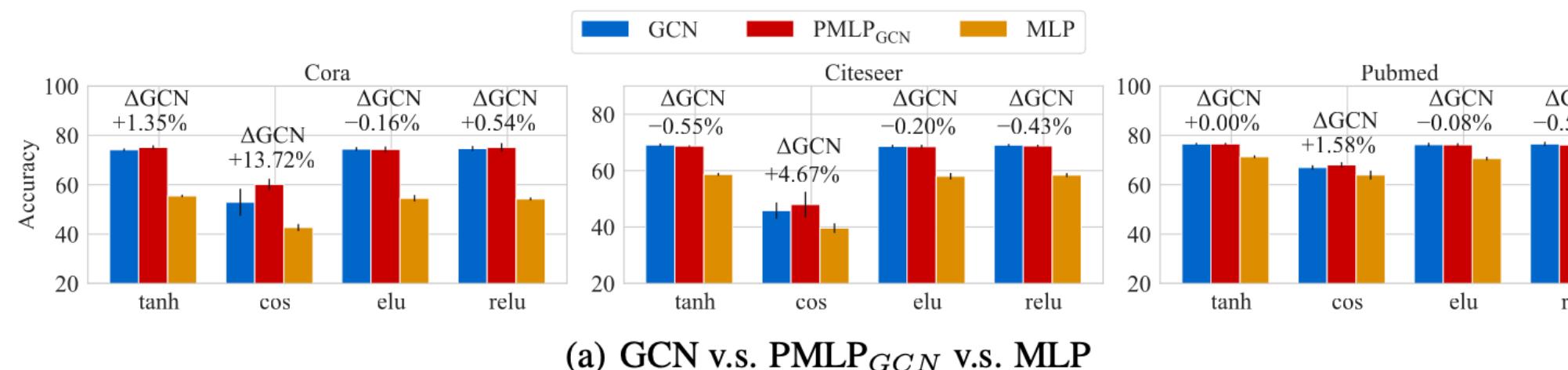
*Q1: What is the impact of model layers and hidden sizes (hyperparameters)?*



Results still hold across different **hyperparameters**.

# Experiments

*Q2/3: What is the impact of different activation functions in FF layers / message passing schemes in MP layers?*



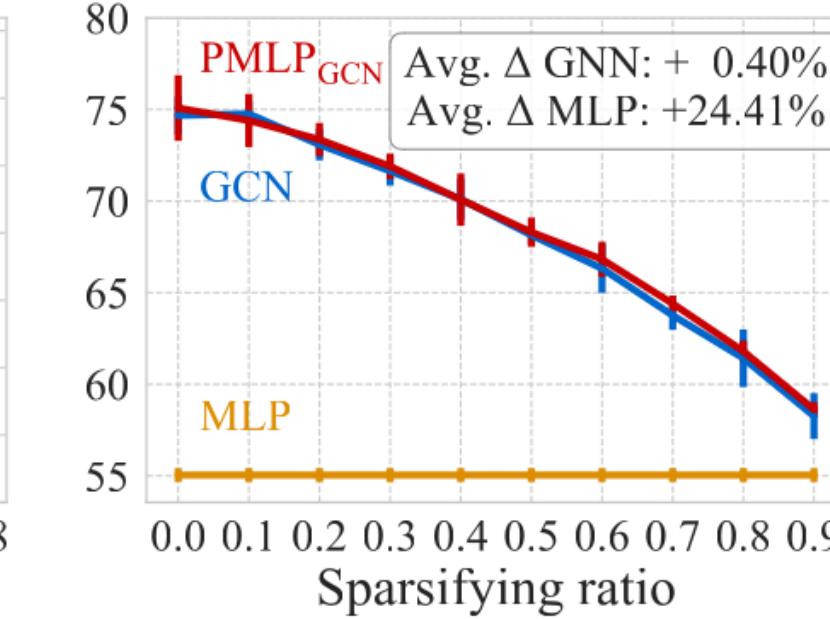
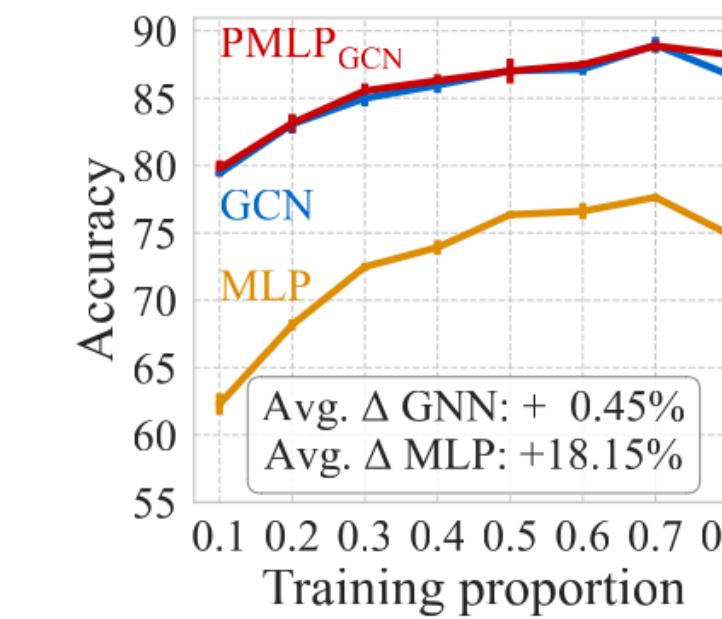
Results still hold across different model implementations.

# Experiments

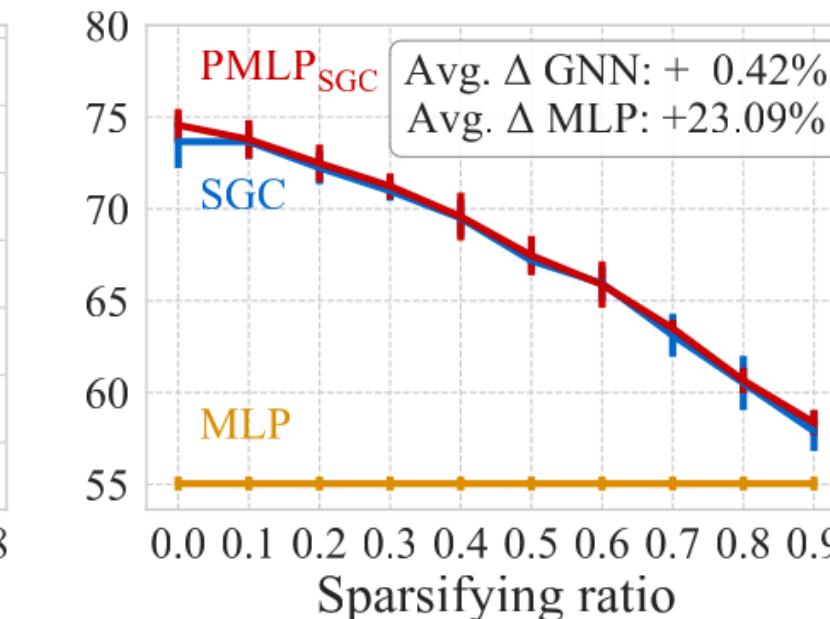
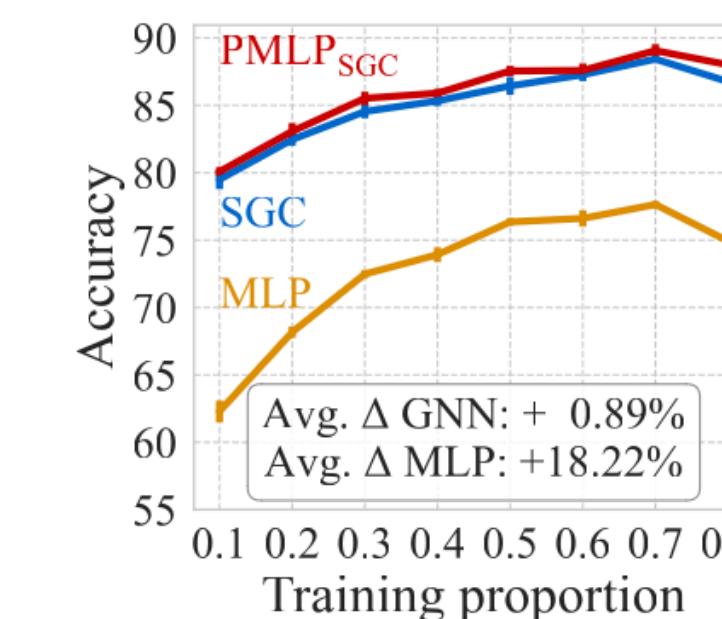
**Q4/5/6: What is the impact of training proportion, graph sparsity and noisy graph structures?**

**Results still hold for different graph properties and PMLP is more robust to graph structure noise.**

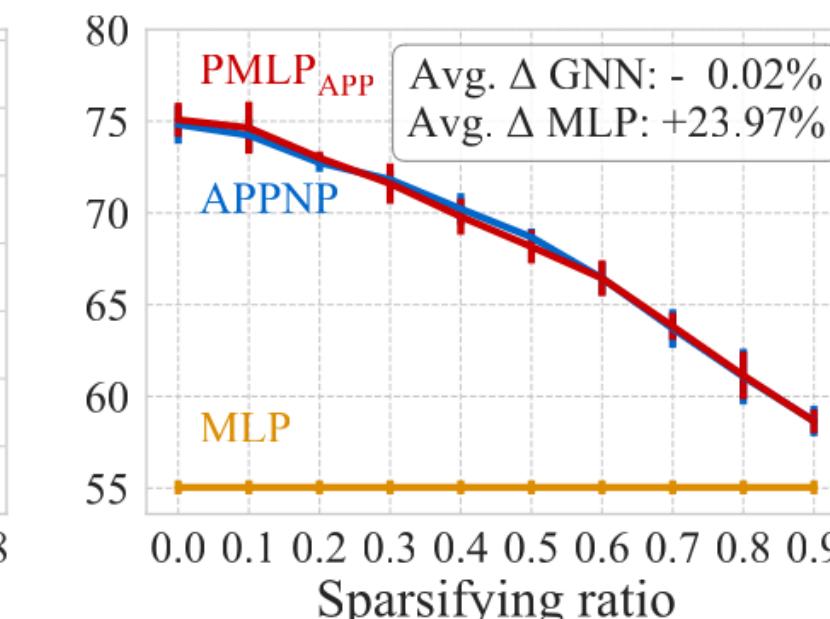
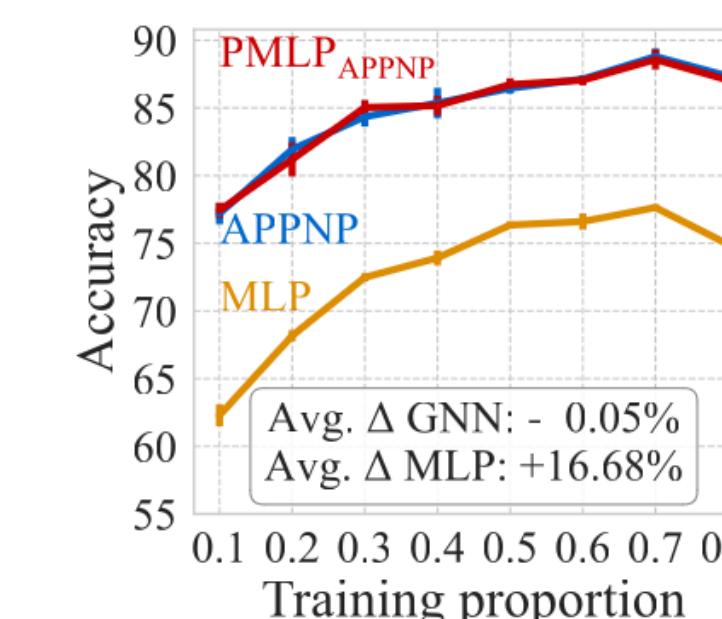
(See more discussions on residual connections, oversmoothing and heterophily in our paper.)



(a) GCN v.s. PMLP<sub>GCN</sub> v.s. MLP



(b) SGC v.s. PMLP<sub>SGC</sub> v.s. MLP



(c) APPNP v.s. PMLP<sub>APP</sub> v.s. MLP

# Quick Guide of Implementation

- PMLP can be implemented in many different ways. The simplest way only requires one line of code.

```
# version A: three models (mlp, pmlp, gnn) in one class
class My_GNN(nn.Module):
    ...
    def forward(self, x, edges, use_conv=True):
        ...
        x = self.feed_forward_layer(x)
        if use_conv:
            x = self.message_passing_layer(x, edges)
        ...
        return x

my_gnn = My_GNN()

# in the training loop
my_gnn.train()
for epoch in range(args.epochs):
    prediction = my_gnn(x, edge, use_conv=False)

# for inference
my_gnn.eval()
prediction = my_gnn(x, edge, use_conv=True)
```

```
# version B: one line of code is all you need
class My_GNN(nn.Module):
    ...
    def forward(self, x, edges):
        ...
        x = self.feed_forward_layer(x)
        if not self.training: # only modified part
            x = self.message_passing_layer(x, edges)
        ...
        return x

my_gnn = My_GNN()

# in the training loop
my_gnn.train()
for epoch in range(args.epochs):
    prediction = my_gnn(x, edges)

# for inference
my_gnn.eval()
prediction = my_gnn(x, edges)
```

# Quick Guide of Implementation

---

- PMLP can be implemented in many different ways. The simplest way only requires one line of code.

*Q: What could we do with PMLP?*

- *Accelerate GNN training.*
- *Empower pre-trained NNs with test-time message passing.*
- *Empower GNN in some scenarios.*
- *Simple and useful tool for analytical purposes.*

# Quick Guide of Implementation

---

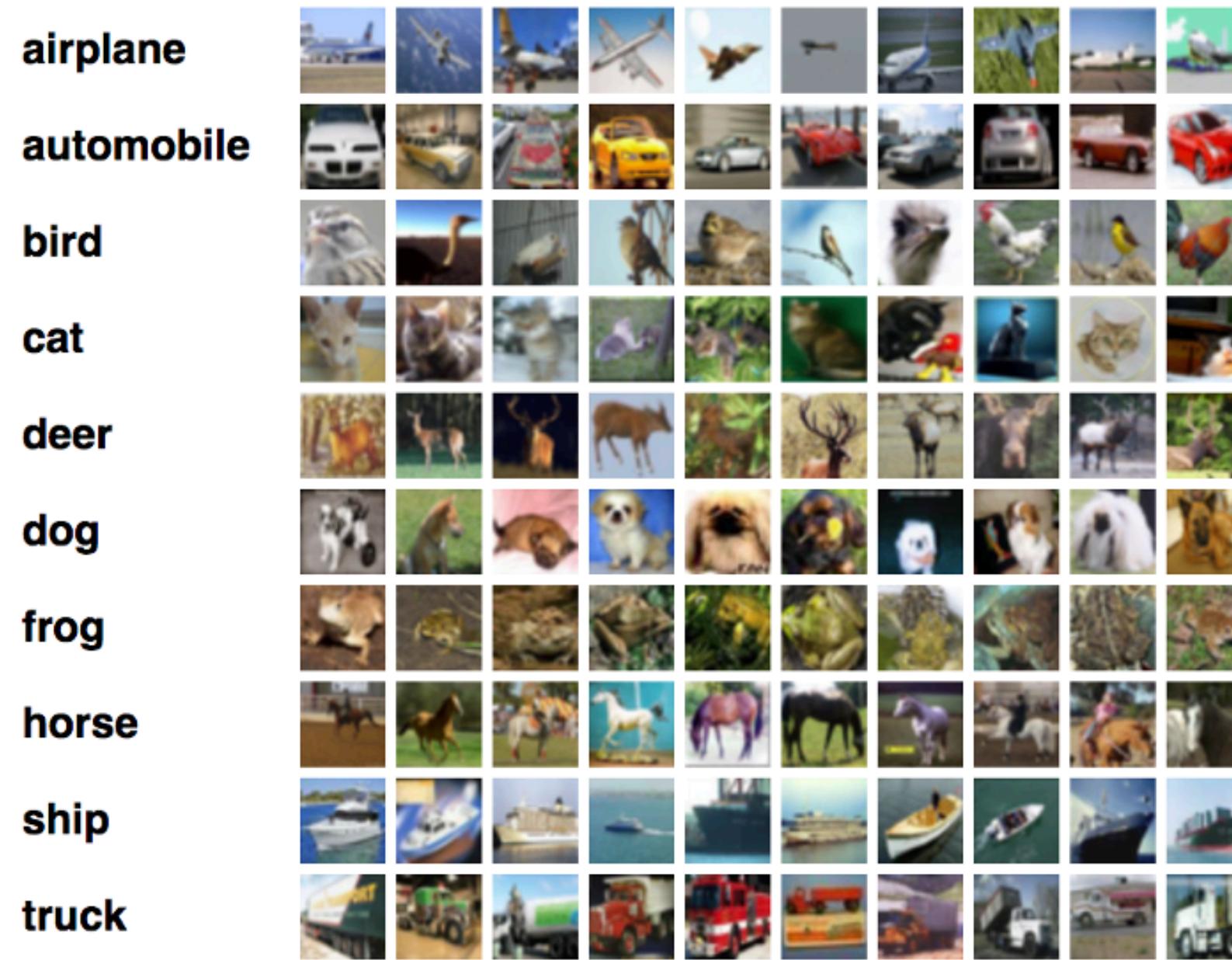
- PMLP can be implemented in many different ways. The simplest way only requires one line of code.

*Q: What are potential extensions of PMLP?*

- *Parameterized MP layers, such as GAT.*
- *Transductive / semi-supervised learning settings.*
- *Other tasks (link prediction, graph classification, etc.)*
- *GNN problems (oversmoothing, heterophily, etc.)*

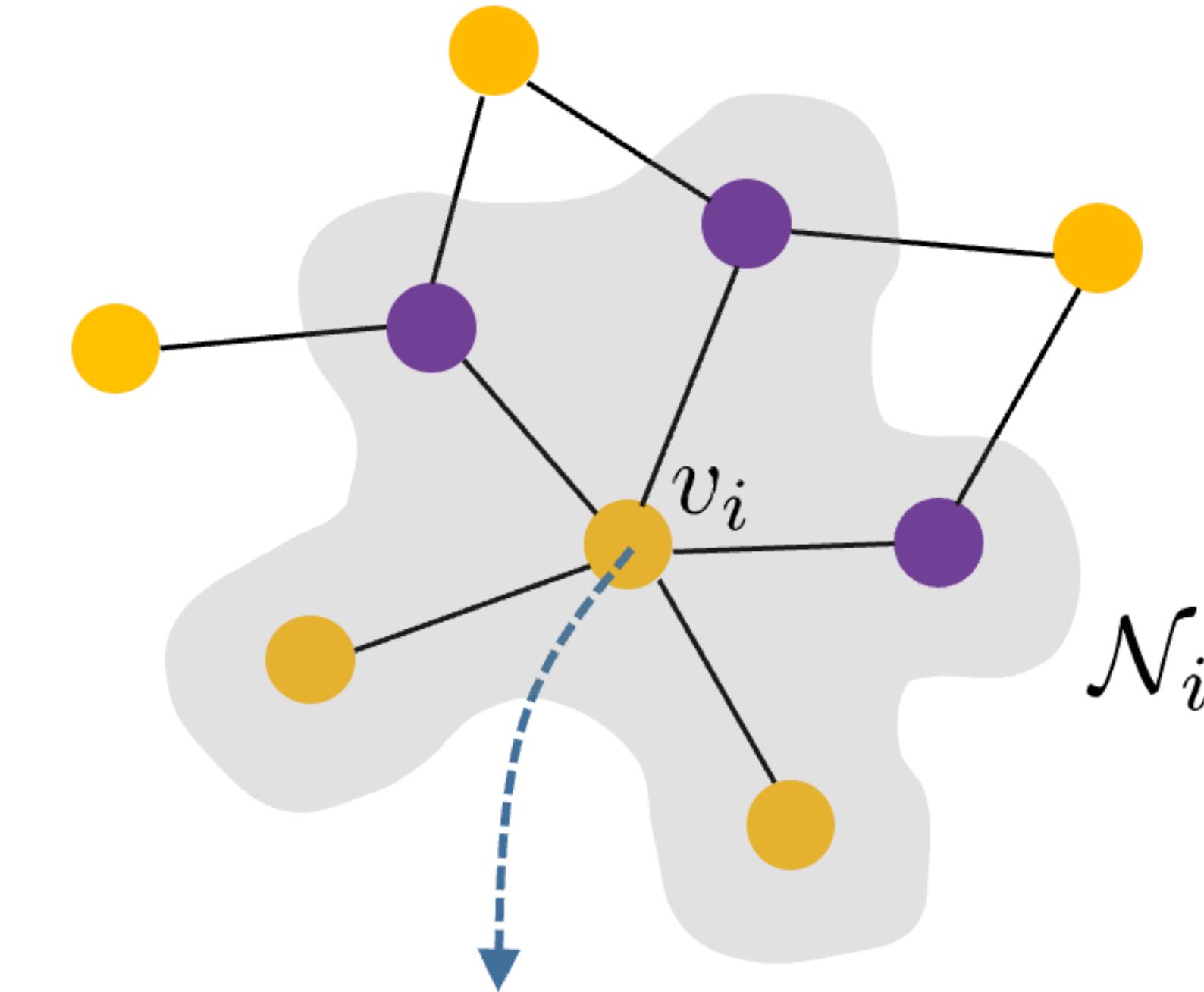
# GNNs are Inherently Good Generalizers

*Why such results are interesting?*



$$(x_i, y_i) \sim p(x, y)$$

each instance is drawn from the same data distribution independently (i.i.d.)



$$(x_i, y_i) \sim p(x, y | \mathcal{N}_i)$$

instances have inter-connection and cannot be treated as i.i.d. samples

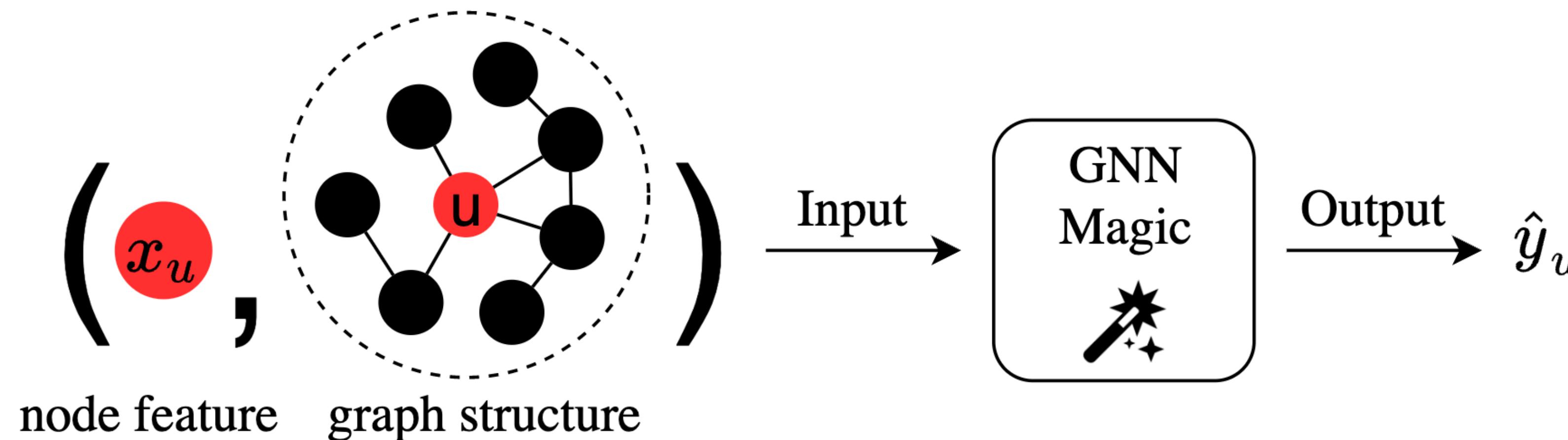
# GNNs are Inherently Good Generalizers

---

*Why such results are interesting?*

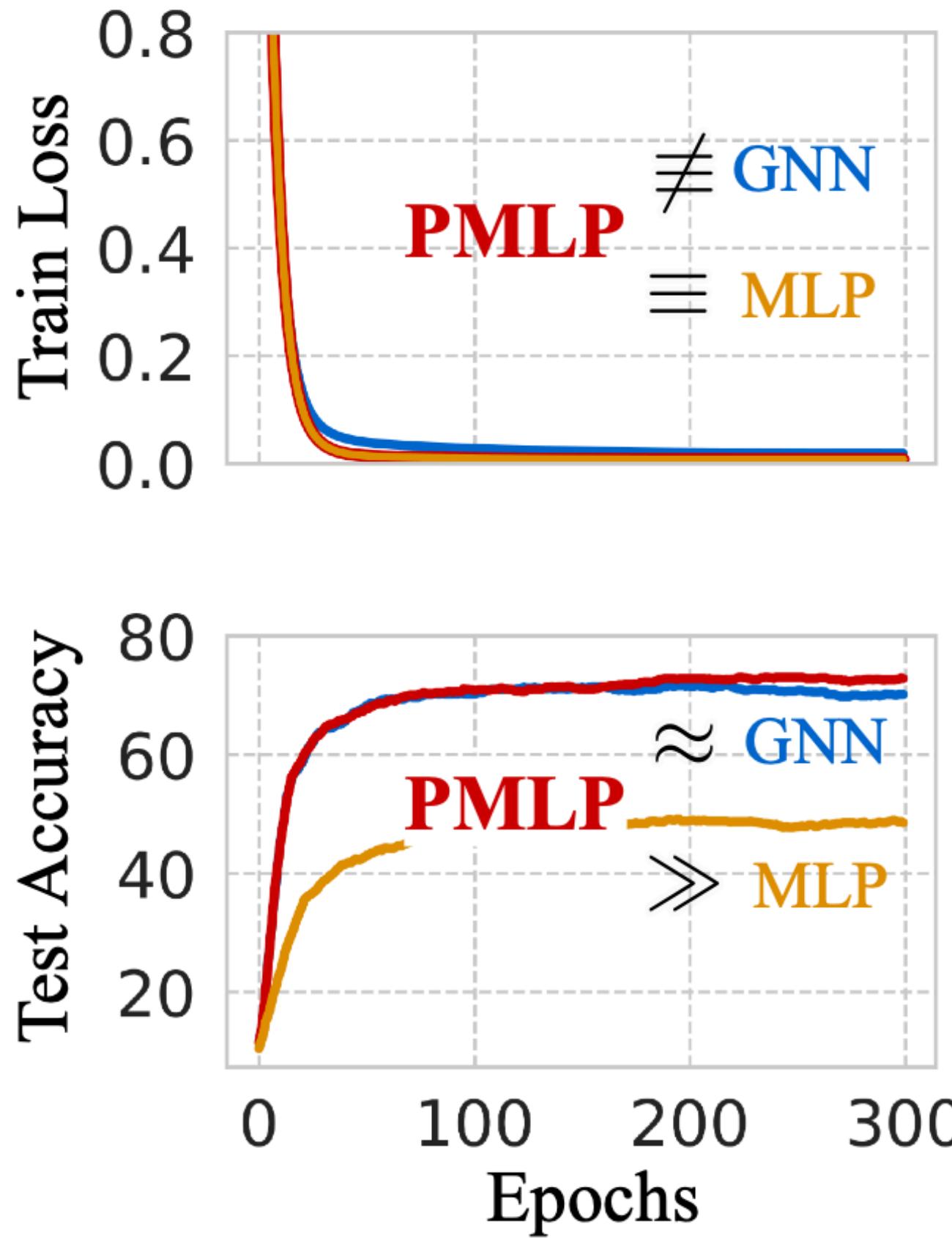
*Conventional wisdom :*

*GNNs are more expressive than MLPs and capable of somehow learning to encode graph information through repeated message passing operations in a way useful for the classification.*



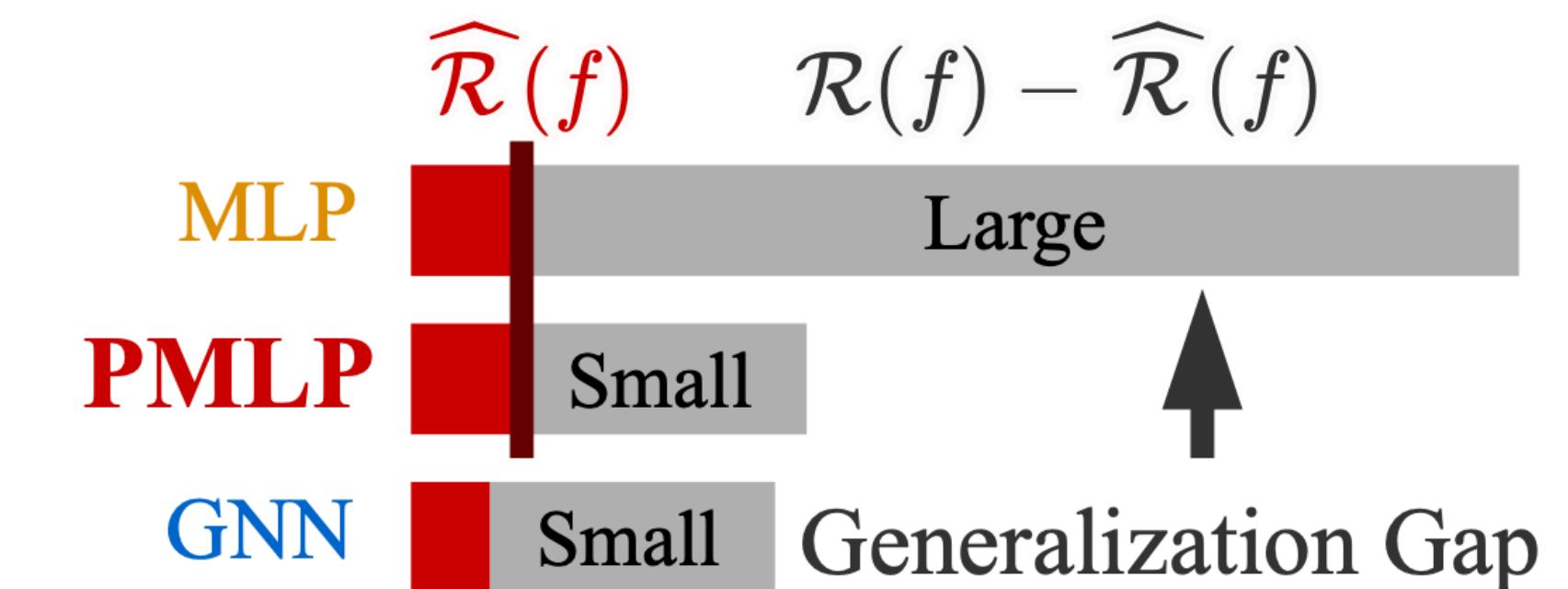
*This might NOT be the key factor for explaining GNNs' success in node-level tasks.*

# GNNs are Inherently Good Generalizers

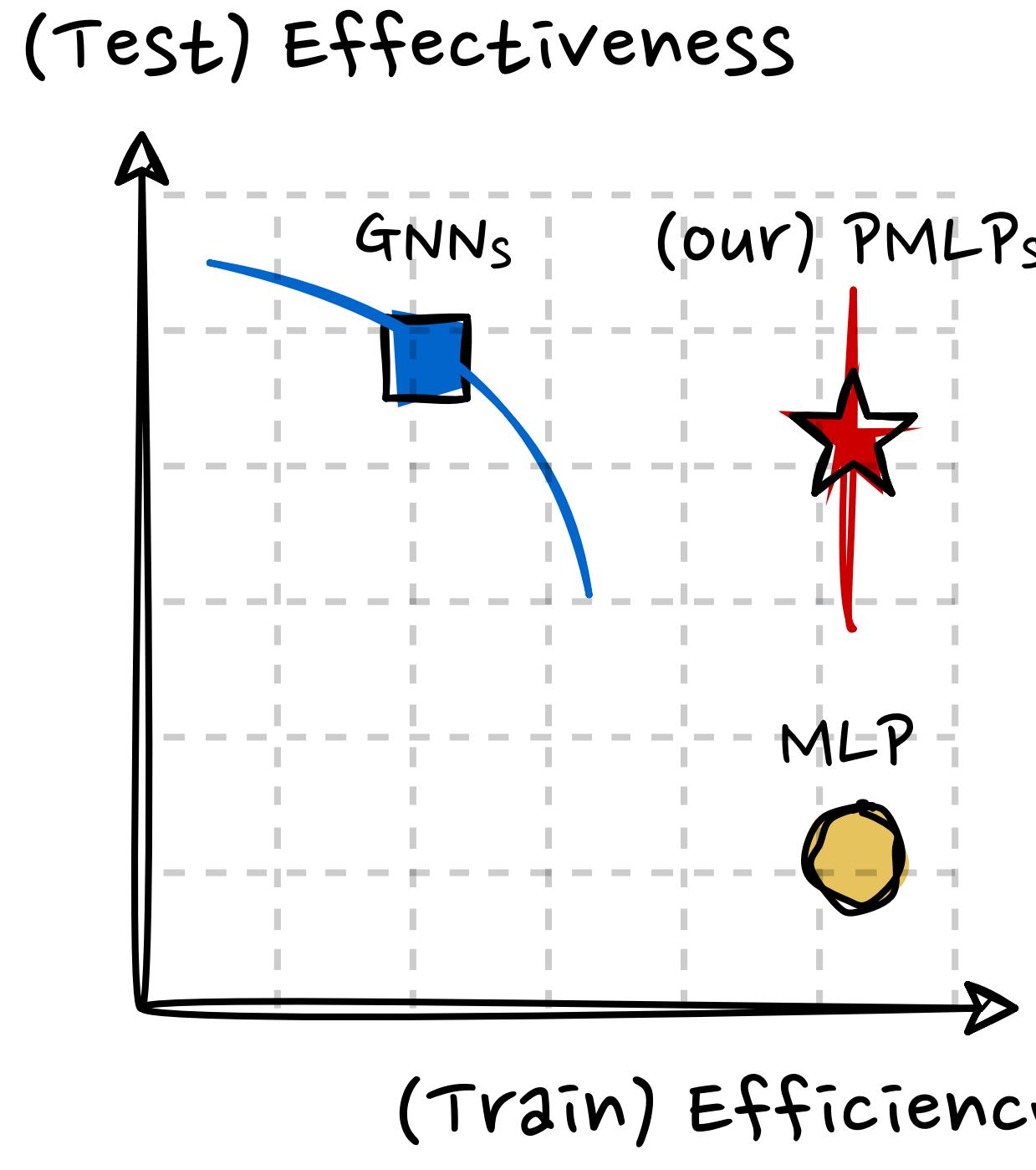


- **Representational Power** : what function class neural networks can approximate and to what extent they can minimize the empirical risk  $\widehat{\mathcal{R}}(f)$ .
- **Generalization** : how well the learned function can generalize to unseen in- and out-of-distribution samples, reflected by the generalization gap  $\mathcal{R}(f) - \widehat{\mathcal{R}}(f)$

*The GNN architecture used in inference inherently encourages lower generalization gap.*



# GNNs are Inherently Good Generalizers



What lies between MLP and GNN?



- Same optimization process and learning dynamics
- Same model expressivity and capability to fit the data
- + Different generalizability due to the GNN architecture in inference  
*(key factor)*



- Same GNN architecture in inference
- + Different optimization process and model expressivity

PMLPs pinpoint the major source of performance gap between GNNs and MLPs in node classification stems from the **inherent generalization capability of GNN architectures**.

# Effects of Model Architectures on Generalization

---

- **(Neural Tangent Kernel) NTK perspective on MLP, PMLP and GNN :** NTK perspective allows us to conveniently study the inherent effects of model architectures due to disentanglement of weights and kernel feature map.

$$f(x; \theta) \approx \underbrace{f(x; \theta_0)}_{\text{random output}} + \underbrace{\nabla_{\theta} f(x; \theta)^{\top} (\theta - \theta_0)}_{\text{feature map}}$$

**learned weights**

As the network width tends to infinity, the feature map becomes constant and is determined by the model architecture (e.g., MLP, GNN and CNN), inducing a kernel called NTK

$$\text{NTK}(x_i, x_j) = \phi_{ntk}(x_i)^{\top} \phi_{ntk}(x_j) = \left\langle \nabla_{\theta} f(x_i; \theta), \nabla_{\theta} f(x_j; \theta) \right\rangle$$

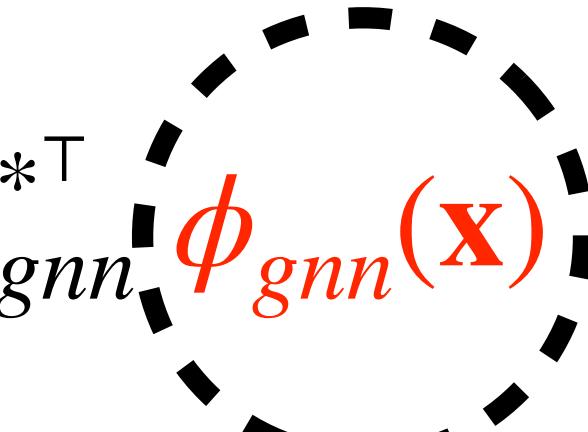
# Effects of Model Architectures on Generalization

---

**Proposition 1.** *MLP and its corresponding PMLP have the same minimum RKHS-norm NTK kernel regression solution, but differ from that of GNNs, i.e.,  $\mathbf{w}_{mlp}^* = \mathbf{w}_{pmlp}^* \neq \mathbf{w}_{gnn}^*$ .*

- **Implications :** PMLP corresponds to transforming the kernel feature map of MLP to that of GNN, while fixing trained MLP weights.

$$f_{mlp}(\mathbf{x}) = \mathbf{w}_{mlp}^{*\top} \phi_{mlp}(\mathbf{x}) \quad f_{pmlp}(\mathbf{x}) = \mathbf{w}_{mlp}^{*\top} \phi_{gnn}(\mathbf{x}) \quad f_{gnn}(\mathbf{x}) = \mathbf{w}_{gnn}^{*\top} \phi_{gnn}(\mathbf{x})$$



The key factor of good testing performance

**Lemma 2.** *The explicit form of GNTK feature map for a two-layer GNN  $\phi_{gnn}(x)$  with average aggregation and ReLU activation in node regression is*

$$\phi_{gnn}(\mathbf{x}_i) = c \sum_{j \in \mathcal{N}_i \cup \{i\}} \left[ \mathbf{X}^\top \mathbf{a}_j \cdot \mathbb{I}_+ \left( \mathbf{w}^{(k)\top} \mathbf{X}^\top \mathbf{a}_j \right), \mathbf{w}^{(k)\top} \mathbf{X}^\top \mathbf{a}_j \cdot \mathbb{I}_+ \left( \mathbf{w}^{(k)\top} \mathbf{X}^\top \mathbf{a}_j \right), \dots \right], \quad (7)$$

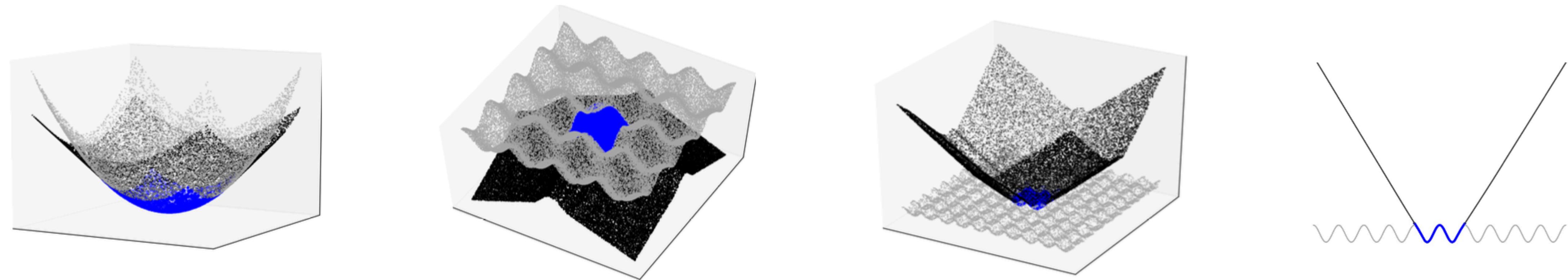
# Effects of Model Architectures on Generalization

---

**Theorem 3.** (*Xu et al., 2021*). Suppose  $f_{mlp}(\mathbf{x})$  is an infinitely-wide two-layer MLP with ReLU trained by square loss. For any direction  $\mathbf{v} \in \mathbb{R}^d$  and step size  $\Delta t > 0$ , let  $\mathbf{x}_0 = t\mathbf{v}$ , we have

$$\left| \frac{(f_{mlp}(\mathbf{x}_0 + \Delta t\mathbf{v}) - f_{mlp}(\mathbf{x}_0)) / \Delta t}{c_{\mathbf{v}}} - 1 \right| = O\left(\frac{1}{t}\right). \quad (8)$$

where  $c_{\mathbf{v}}$  is a constant linear coefficient. That is, as  $t \rightarrow \infty$ ,  $f_{mlp}(\mathbf{x}_0)$  converges to a linear function.



# Effects of Model Architectures on Generalization

---

**Theorem 4.** Suppose  $f_{pmlp}(\mathbf{x})$  is an infinitely-wide two-layer MLP with ReLU activation trained using squared loss, and adds average message passing layer before each feed-forward layer in the testing phase. For any direction  $\mathbf{v} \in \mathbb{R}^d$  and step size  $\Delta t > 0$ , let  $\mathbf{x}_0 = t\mathbf{v}$ , and as  $t \rightarrow \infty$ , we have

$$(f_{pmlp}(\mathbf{x}_0 + \Delta t\mathbf{v}) - f_{pmlp}(\mathbf{x}_0)) / \Delta t \rightarrow c_{\mathbf{v}} \sum_{i \in \mathcal{N}_0 \cup \{0\}} (\tilde{d} \cdot \tilde{d}_i)^{-1}. \quad (9)$$

where  $c_{\mathbf{v}}$  is the same constant as in Theorem 3,  $\tilde{d}_0 = \tilde{d}$  is the node degree (with self-connection) of  $\mathbf{x}_0$ , and  $\tilde{d}_i$  is the node degree of its neighbors.

- **Interpretation:** Similar to MLP, both PMLP and GNN (with sum/mean pooling) converge to linear functions when testing samples are far away from the training data.

# Effects of Model Architectures on Generalization

---

**Theorem 5.** Suppose all node features are normalized, and the cosine similarity of node  $\mathbf{x}_i$  and the average of its neighbors is denoted as  $\alpha_i \in [0, 1]$ . Then, the convergence rate for  $f_{pmlp}(\mathbf{x})$  is

$$\left| \frac{(f_{pmlp}(\mathbf{x}_0 + \Delta t \mathbf{v}) - f_{pmlp}(\mathbf{x}_0)) / \Delta t}{c_{\mathbf{v}} \sum_{i \in \mathcal{N}_0 \cup \{0\}} (\tilde{d} \cdot \tilde{d}_i)^{-1}} - 1 \right| = O \left( \frac{1 + (\tilde{d}_{max} - 1) \sqrt{1 - \alpha_{min}^2}}{t} \right). \quad (10)$$

where  $\alpha_{min} = \min\{\alpha_i\}_{i \in \mathcal{N}_0 \cup \{0\}} \in [0, 1]$ , and  $\tilde{d}_{max} \geq 1$  denotes the maximum node degree in the testing node  $\mathbf{x}_0$ 's neighbors (including itself).

**Theorem 3. (Xu et al., 2021).** Suppose  $f_{mlp}(\mathbf{x})$  is an infinitely-wide two-layer MLP with ReLU trained by square loss. For any direction  $\mathbf{v} \in \mathbb{R}^d$  and step size  $\Delta t > 0$ , let  $\mathbf{x}_0 = t\mathbf{v}$ , we have

$$\left| \frac{(f_{mlp}(\mathbf{x}_0 + \Delta t \mathbf{v}) - f_{mlp}(\mathbf{x}_0)) / \Delta t}{c_{\mathbf{v}}} - 1 \right| = O\left(\frac{1}{t}\right). \quad (8)$$

where  $c_{\mathbf{v}}$  is a constant linear coefficient. That is, as  $t \rightarrow \infty$ ,  $f_{mlp}(\mathbf{x}_0)$  converges to a linear function.

# Effects of Model Architectures on Generalization

---

**Theorem 5.** Suppose all node features are normalized, and the cosine similarity of node  $\mathbf{x}_i$  and the average of its neighbors is denoted as  $\alpha_i \in [0, 1]$ . Then, the convergence rate for  $f_{pmlp}(\mathbf{x})$  is

$$\left| \frac{(f_{pmlp}(\mathbf{x}_0 + \Delta t \mathbf{v}) - f_{pmlp}(\mathbf{x}_0)) / \Delta t}{c_{\mathbf{v}} \sum_{i \in \mathcal{N}_0 \cup \{0\}} (\tilde{d} \cdot \tilde{d}_i)^{-1}} - 1 \right| = O \left( \frac{1 + (\tilde{d}_{max} - 1) \sqrt{1 - \alpha_{min}^2}}{t} \right). \quad (10)$$

where  $\alpha_{min} = \min\{\alpha_i\}_{i \in \mathcal{N}_0 \cup \{0\}} \in [0, 1]$ , and  $\tilde{d}_{max} \geq 1$  denotes the maximum node degree in the testing node  $\mathbf{x}_0$ 's neighbors (including itself).

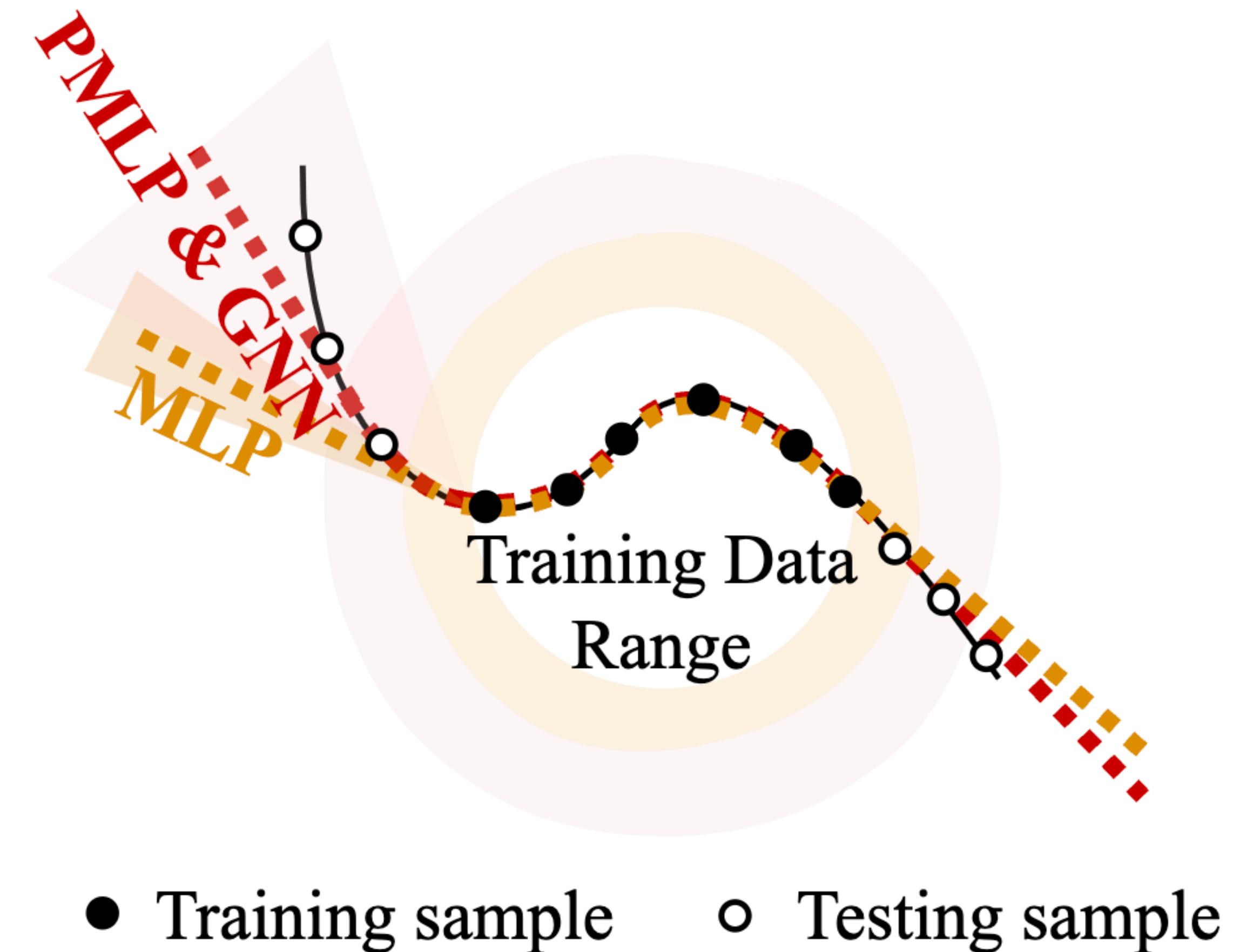
- **Interpretation:** PMLP and GNN's convergence rates (to linear models) are smaller than that of MLP. This indicates they are less vulnerable to linearization, and prone to generalize better on testing samples near the training data.

# Effects of Model Architectures on Generalization

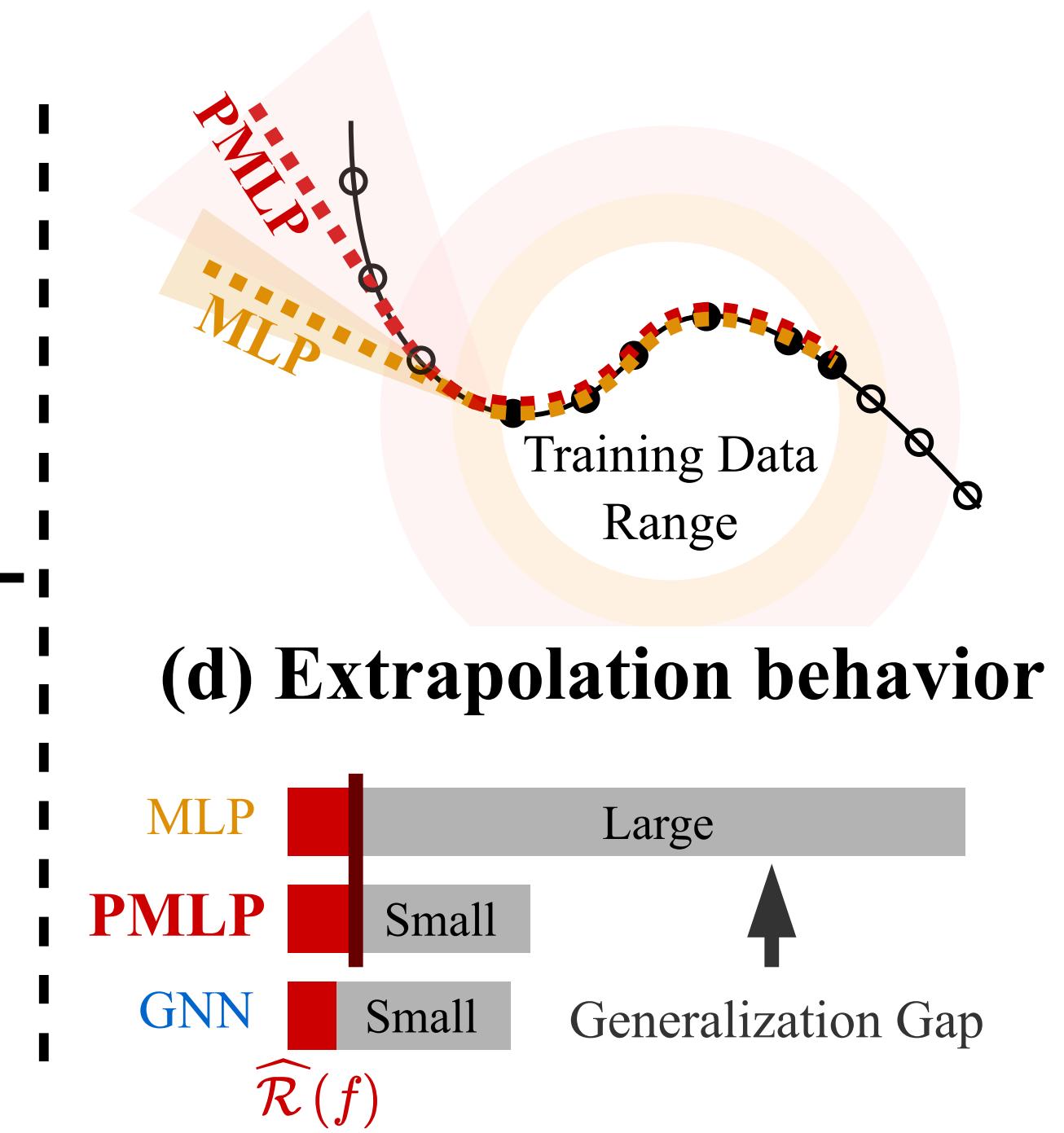
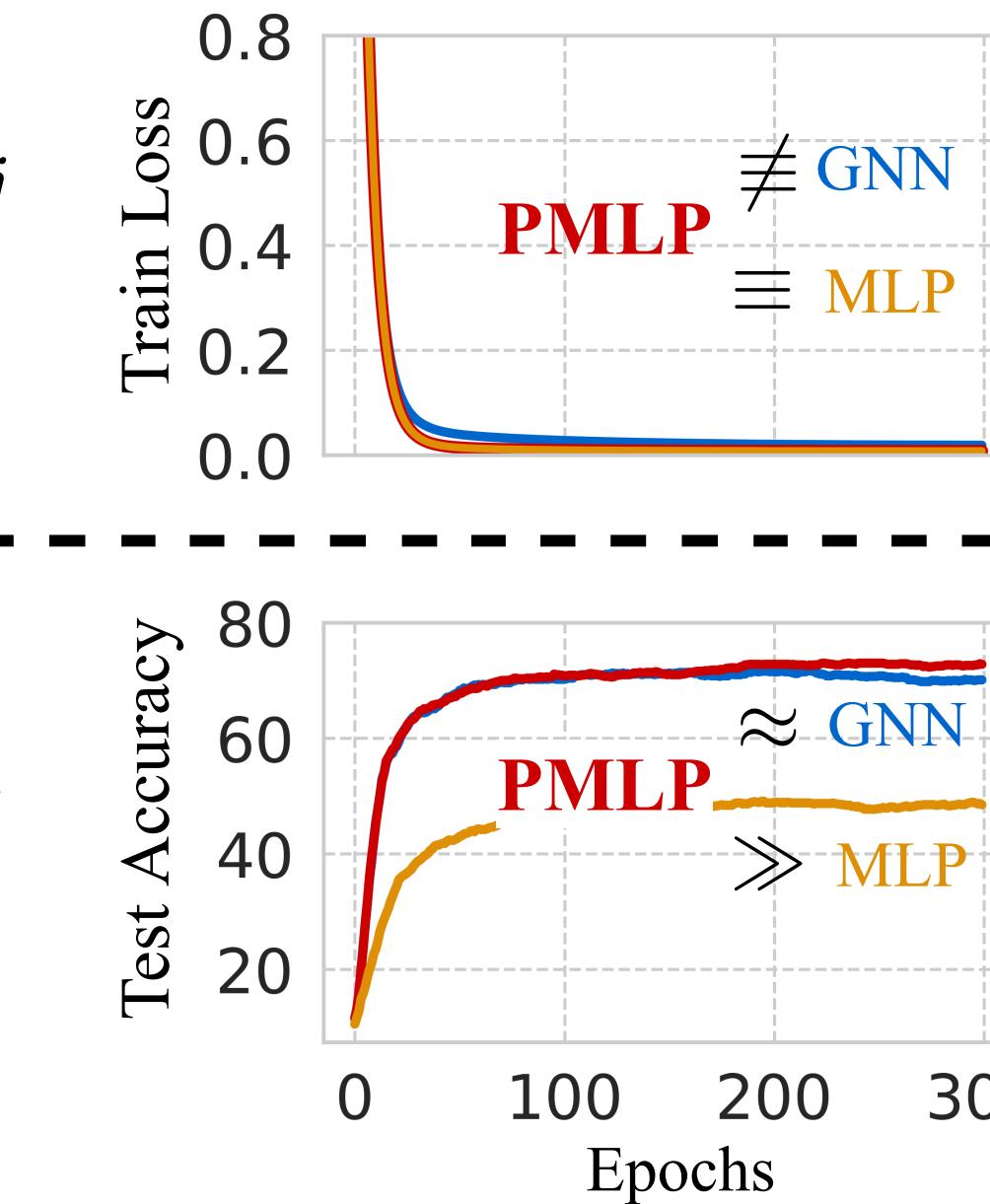
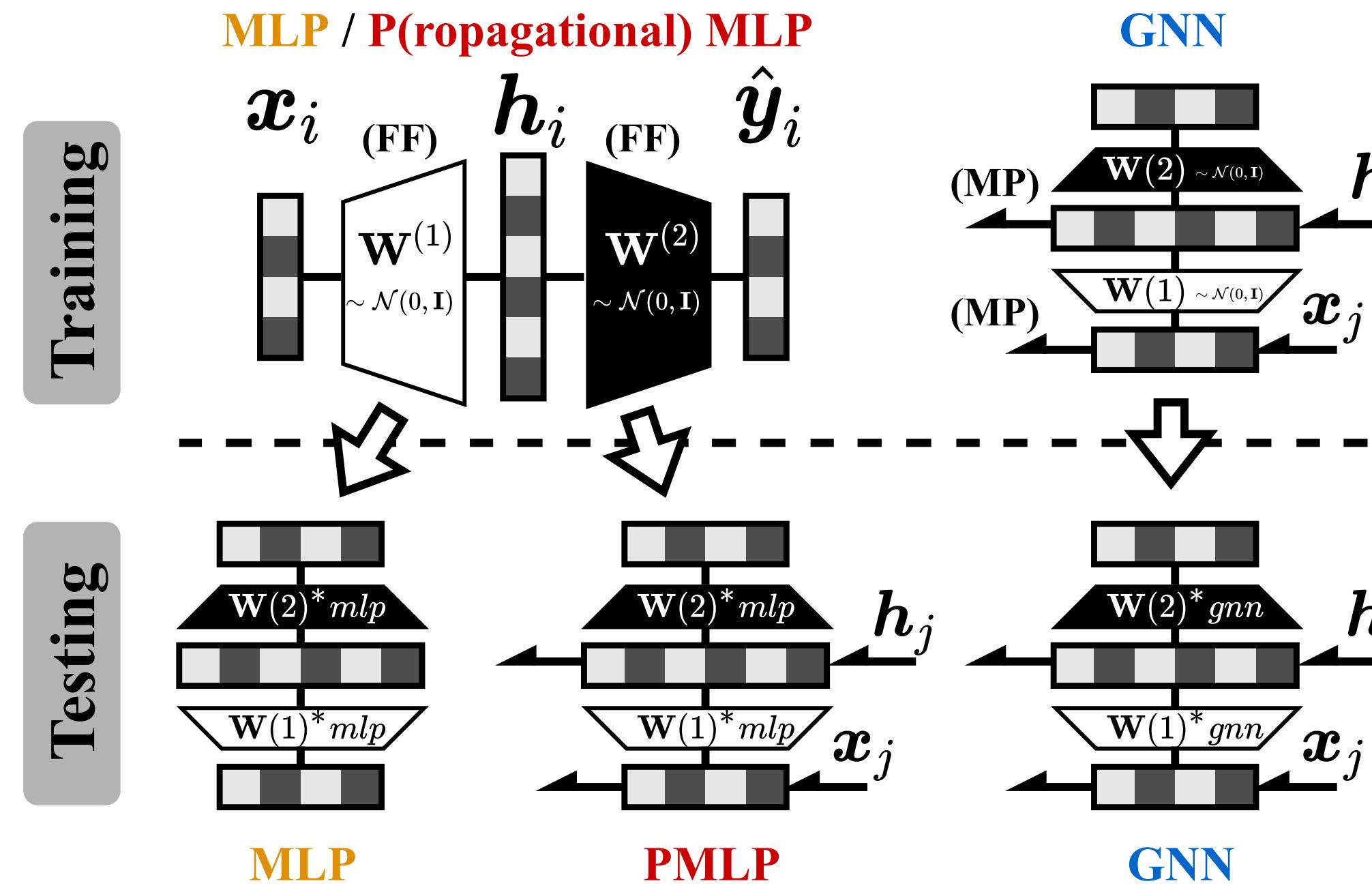
---

- **Theoretical Insight (informal)**

Due to their architectures used in inference, both GNN and PMLP can better extrapolate out-of-distribution testing nodes for node-level tasks.



# Thank you !



*Code and related materials in <https://github.com/chr26195/PMLP>*