

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

**Chenxiao Yang, Qitian Wu, Jiahua Wang, Junchi Yan**

**Department of Computer Science and Engineering, Shanghai Jiao Tong University**



# MLP v.s. GNN : Architectural Connection

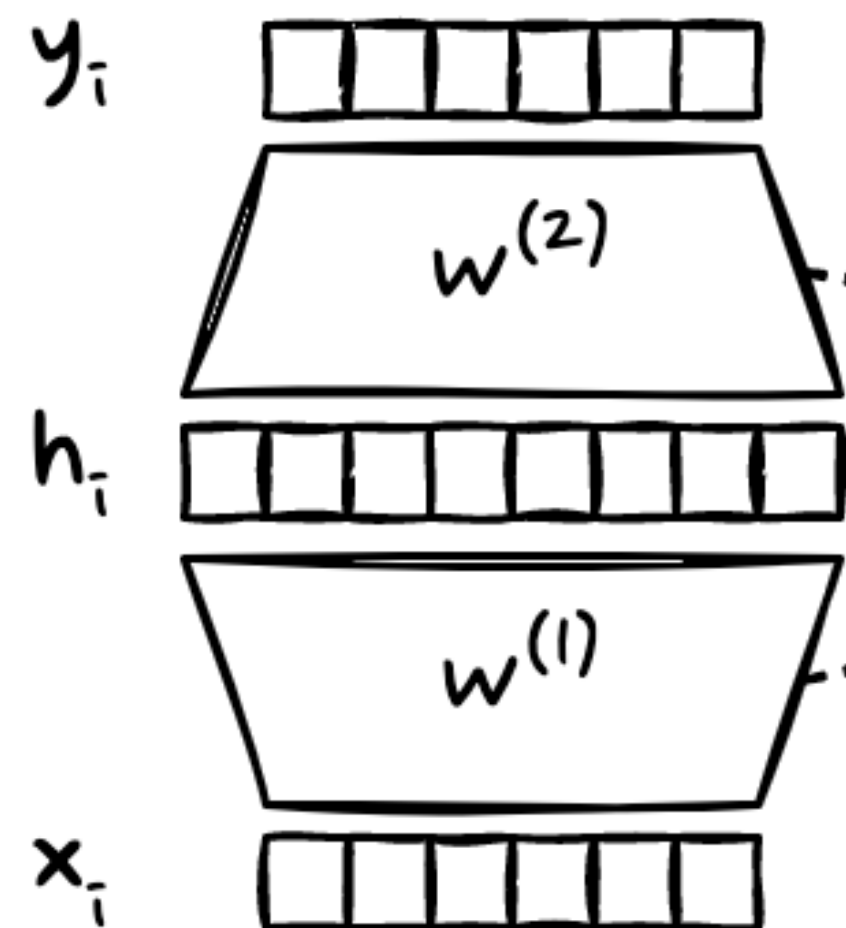
**GNN Layer :** (MP):  $\tilde{\mathbf{h}}_u^{(l-1)} = \sum_{v \in \mathcal{N}_u \cup \{u\}} a_G(u, v) \cdot \mathbf{h}_u^{(l-1)},$

**MLP Layer :** ~~(MP):  $\tilde{\mathbf{h}}_u^{(l-1)} = \sum_{v \in \mathcal{N}_u \cup \{u\}} a_G(u, v) \cdot \mathbf{h}_u^{(l-1)}$~~

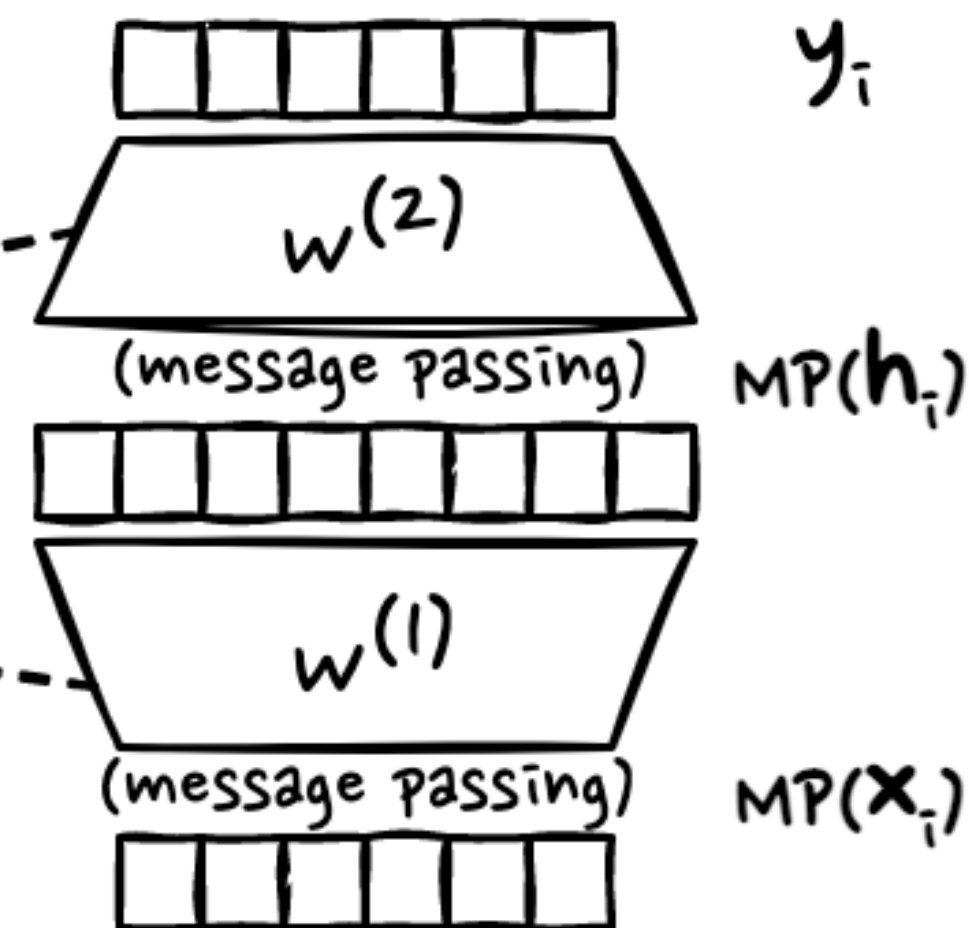
(FF):  $\mathbf{h}_u^{(l)} = \psi^{(l)} \left( \tilde{\mathbf{h}}_u^{(l-1)} \right)$

(FF):  $\mathbf{h}_u^{(l)} = \psi^{(l)} \left( \mathbf{h}_u^{(l-1)} \right)$

Multi-Layer Perceptrons



Graph Neural Networks  
(GCN-style)

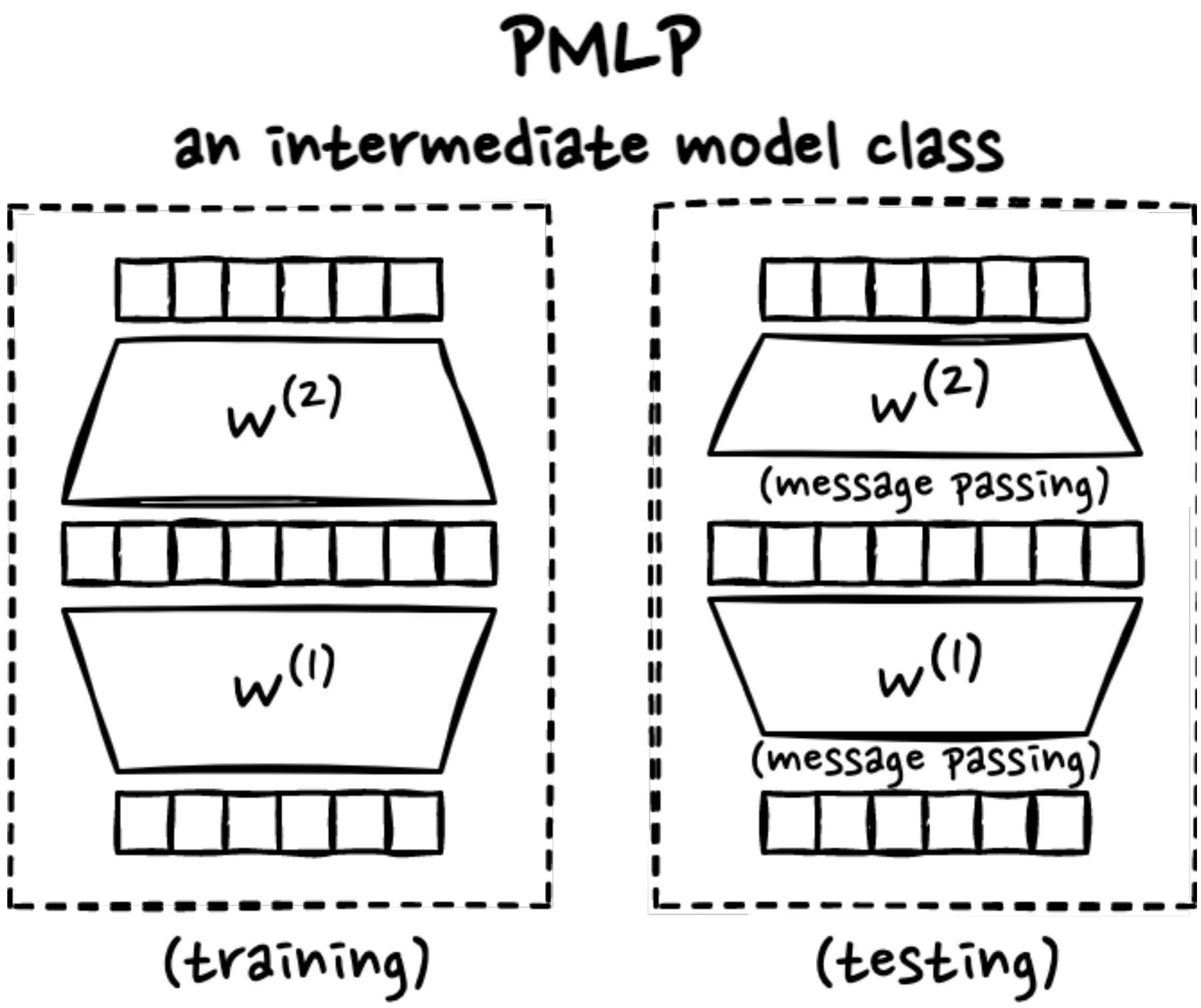


Same  
Parameter  
Space

## Observation

GNN and MLP share the same feed-forward backbone and parameter (sub)space.

# MLP - (?) - GNN : Propagational MLP (PMLP)



**PMLP** : Training with **MLP**  
Inference with **GNN**

## Instantiations of PMLP

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>APNP</sub>		APNP: $\text{Multi-MP}(\psi(\{\mathbf{x}_v\}_{v \in \mathcal{V}}))$
PMLP <sub>GCNII</sub>	ResNet	GCNII
PMLP <sub>JKNet</sub>	MLP+JK	JKNet

⋮

(could be extended to other GNN architectures  
with some modifications)

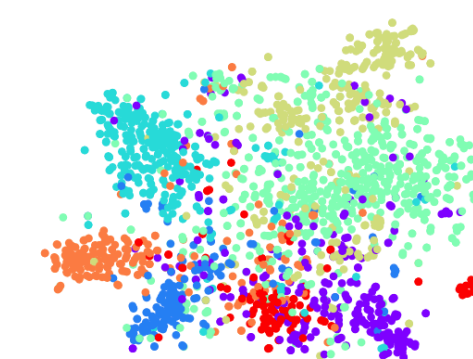
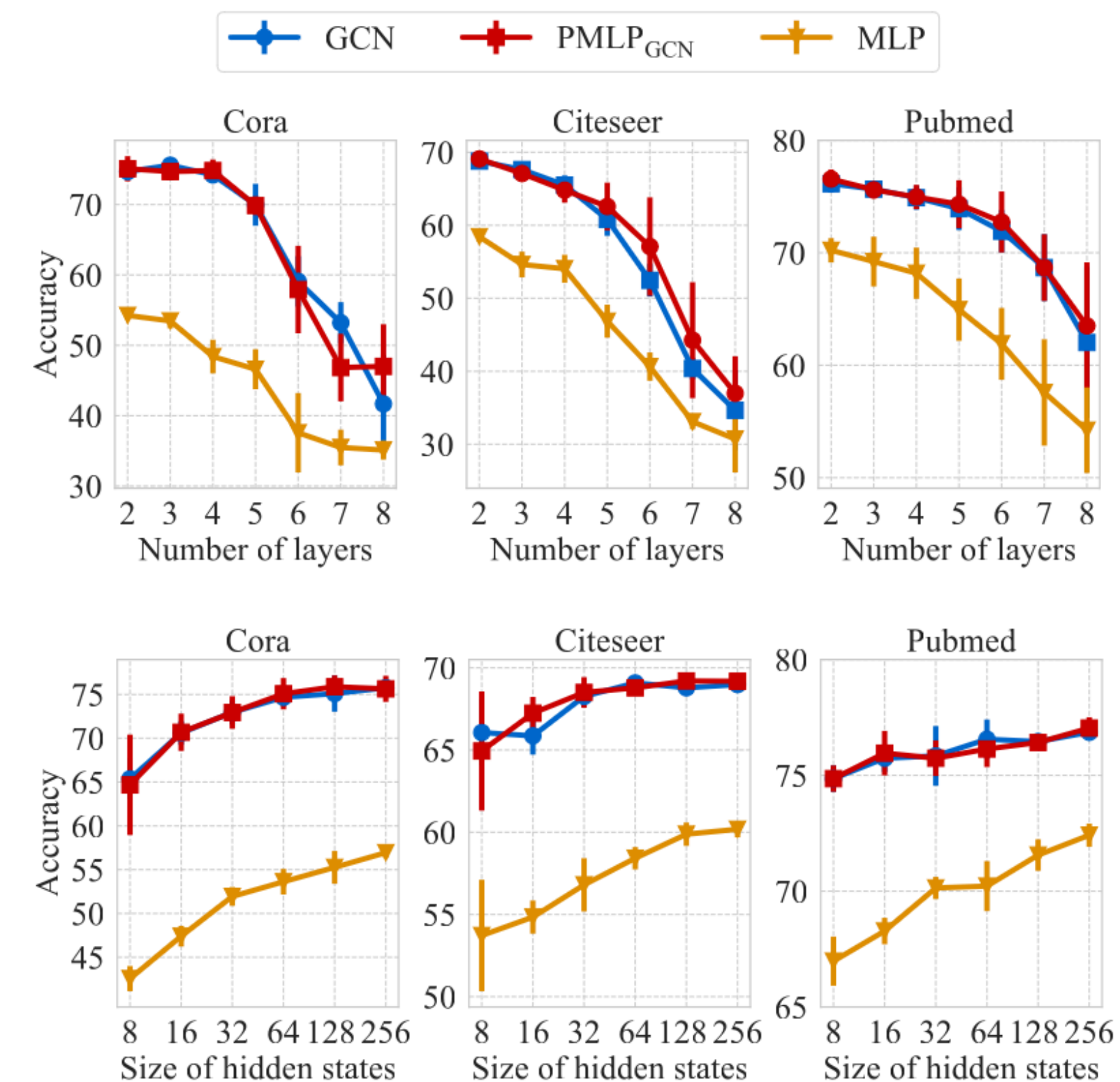
⋮



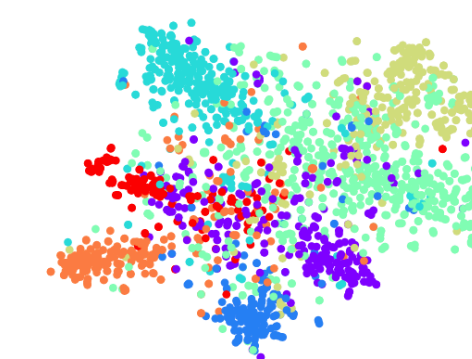
# Empirical Evaluation

Dataset #Nodes	Cora 2,708	Citeseer 3,327	Pubmed 19,717	A-Photo 7,650	A-Computer 13,752	Coauthor-CS 18,333	Coauthor-Physics 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
	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
	APNP	75.02 ± 2.17	66.58 ± 0.77	76.48 ± 0.49	89.51 ± 0.86	78.29 ± 0.55	91.64 ± 0.34
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
	<b>PMLP<sub>GCN</sub></b>	<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>
	$\Delta_{GNN}$	<b>+1.39%</b>	<b>+0.59%</b>	<b>-0.65%</b>	<b>-0.66%</b>	<b>-0.94%</b>	<b>-0.03%</b>
	$\Delta_{MLP}$	<b>+37.18%</b>	<b>+21.00%</b>	<b>+7.49%</b>	<b>+17.84%</b>	<b>+23.75%</b>	<b>+4.86%</b>
	<b>PMLP<sub>SGC</sub></b>	<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>
	$\Delta_{GNN}$	<b>+1.46%</b>	<b>+0.48%</b>	<b>+0.03%</b>	<b>+3.69%</b>	<b>-3.09%</b>	<b>-0.16%</b>
	$\Delta_{MLP}$	<b>+35.70%</b>	<b>+20.39%</b>	<b>+7.43%</b>	<b>+14.40%</b>	<b>+18.47%</b>	<b>+4.09%</b>
	<b>PMLP<sub>APP</sub></b>	<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>
	$\Delta_{GNN}$	<b>+1.09%</b>	<b>+1.41%</b>	<b>-0.24%</b>	<b>-1.16%</b>	<b>-0.28%</b>	<b>+0.00%</b>
	$\Delta_{MLP}$	<b>+37.14%</b>	<b>+20.14%</b>	<b>+7.83%</b>	<b>+17.01%</b>	<b>+23.78%</b>	<b>+4.72%</b>

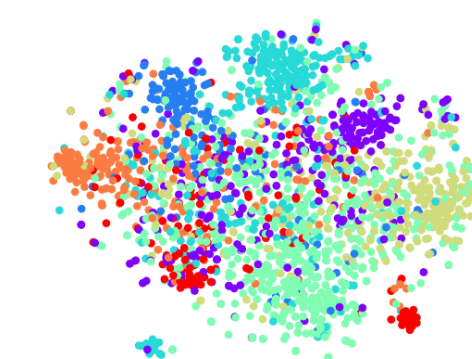
1. PMLPs significantly outperform MLP.
2. PMLPs perform on par with or exceed GNNs in inductive learning setting.
3. PMLPs speed up training up to 65 times.
4. PMLPs are more robust to graph structural noises.



(a) GCN



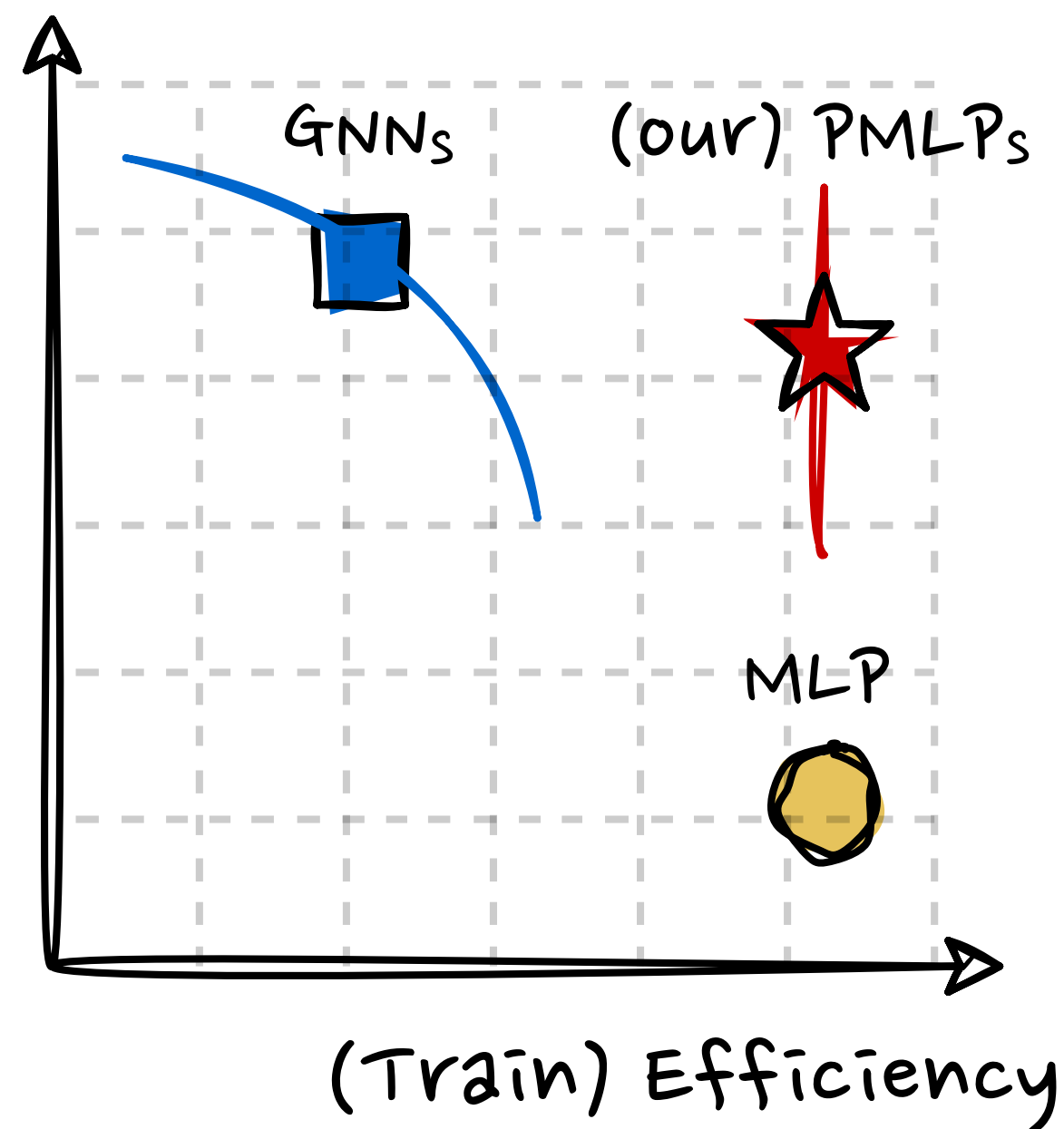
(b) PMLP<sub>GCN</sub>





(c) MLP

# GNNs are Inherently Good Generalizers

(Test) Effectiveness



What lies between MLP and GNN?

MLP   $\xrightarrow{\text{large gap}}$  PMLP  :

- 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)

PMLP   $\xrightarrow{\text{small gap}}$  GNN  :

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

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

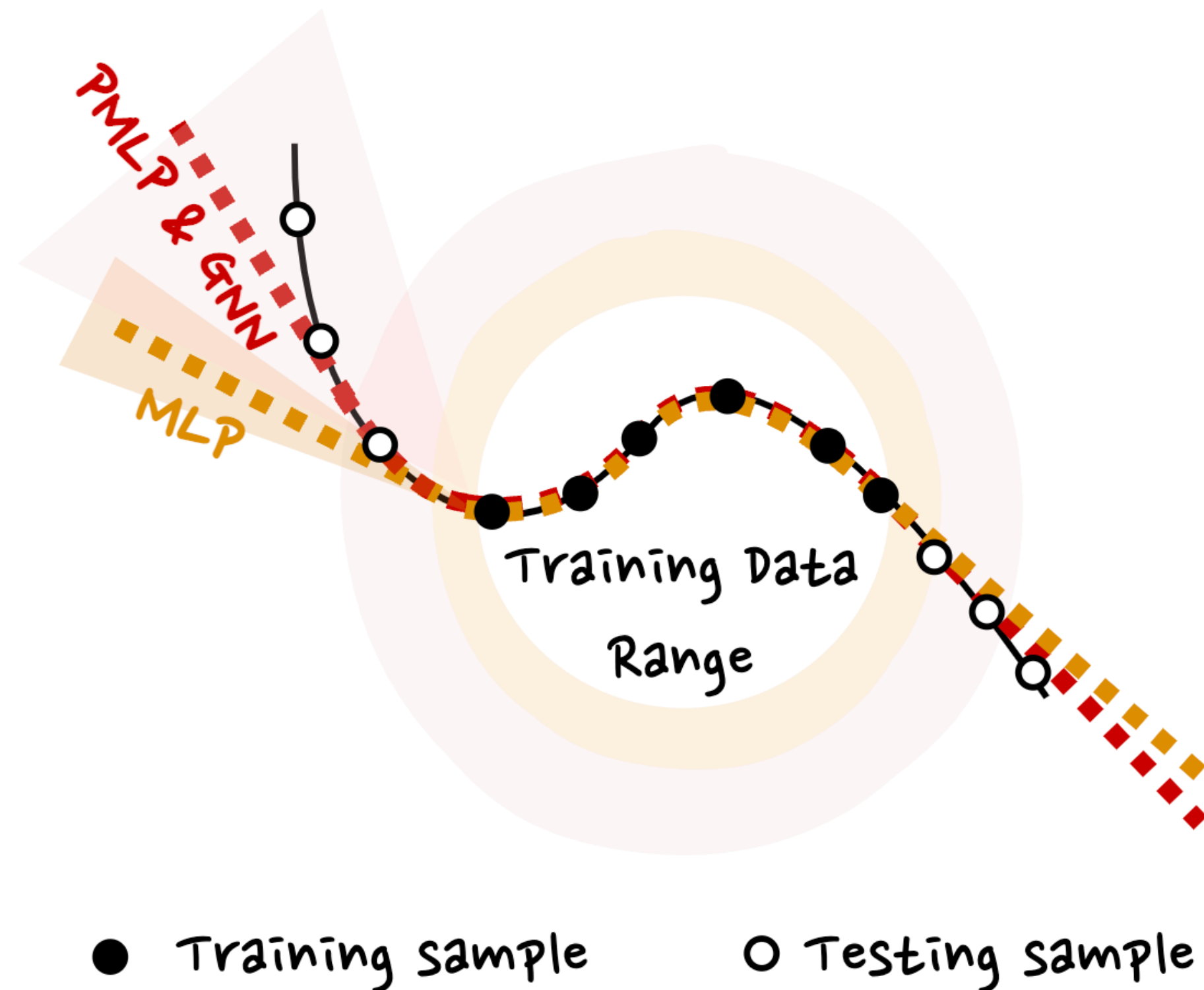
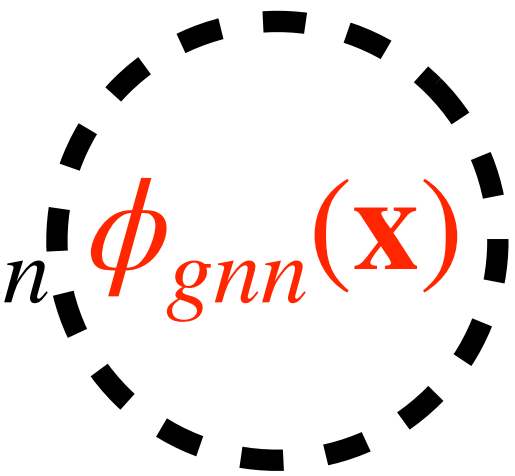


# Effects of Model Architectures in Extrapolation

## Comparison in NTK perspective

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

The key factor



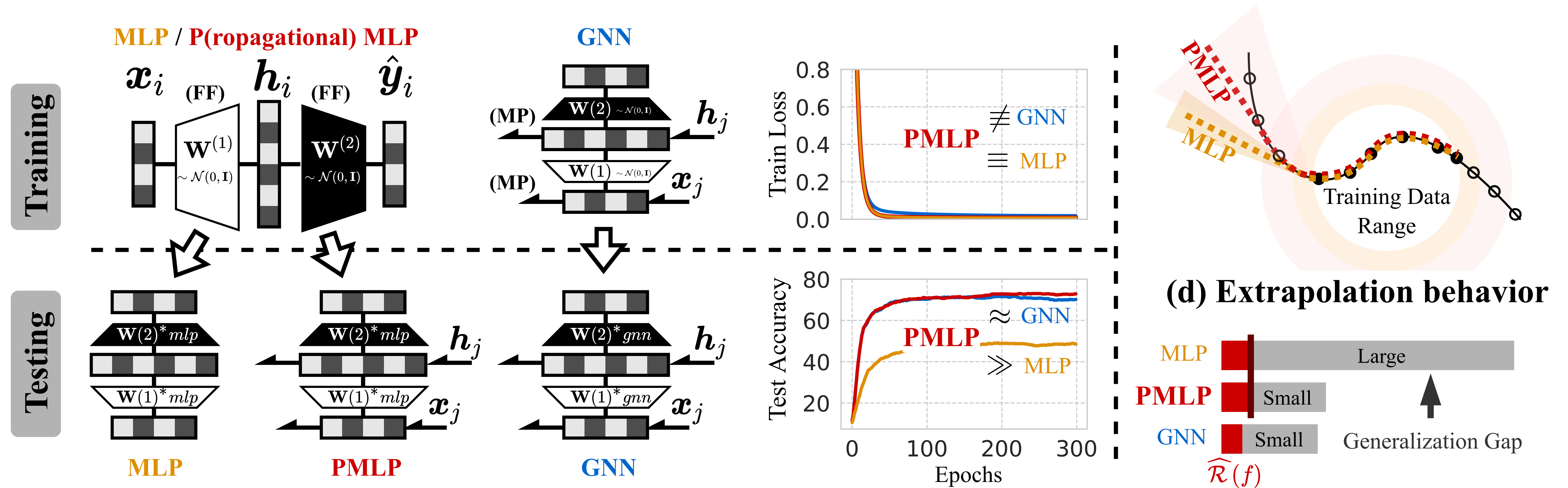
## Theorem 4

Alike MLP, both PMLP and GNN (in GCN-style) eventually converge to linear functions when testing samples are far away from the training data support.

## Theorem 5

PMLP and GNN's convergence rates (to linear models) are smaller than MLP due to message passing at each layer. This indicates they are less more vulnerable to linearization, and prone to generalize to testing samples near the training data.

# Summary and Takeaways



(a) Illustration of model architecture

(b) Learning behavior

(c) Intrinsic generalizability of GNN

1. Vanilla MLP with post-training message passing can be as effective as GNN.
2. GNN inherently extrapolates better than MLP because of its architecture.
3. It is important to understand GNN generalization for understanding its success.