# On Provable Benefits of Depth in Training Graph Convolutional Networks

**Weilin Cong**
Pennsylvania State University
wxc272@psu.edu

**Morteza Ramezani**
Pennsylvania State University
morteza@cse.psu.edu

**Mehrdad Mahdavi**
Pennsylvania State University
mzm616@psu.edu

# Motivation

- Graph neural networks have achieved state-of-the-art performance in many graph-structured applications.

- Existing GNNs are limited to very shallow structures because GNNs suffer from performance degradation issue as the number of layers increases.

- The conventional wisdom is that adding the number of layers cause **over-smoothing**.

- We observe that *there exists a discrepancy between the theoretical understanding of the inherent capabilities of GNN and their practical performance.*

# Motivation

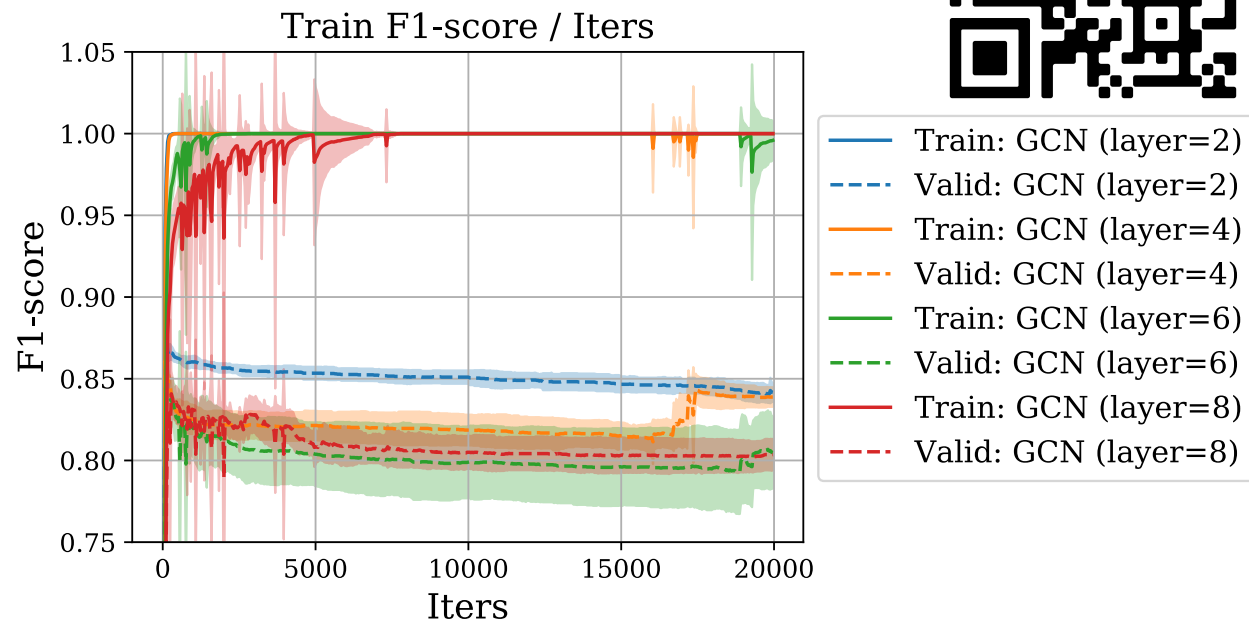- Experiment observations

```python
import dgl.data

dataset = dgl.data.CoraGraphDataset()
print('Number of categories:', dataset.num_classes)
```

```python
from dgl.nn import GraphConv

class GCN(nn.Module):
    def __init__(self, in_feats, h_feats, num_classes, num_layers=2):
        super(GCN, self).__init__()
        self.convs = nn.ModuleList()
        self.num_layers = num_layers

        self.convs.append(GraphConv(in_feats, h_feats))
        for _ in range(num_layers-2):
            self.convs.append(GraphConv(h_feats, h_feats))
        self.convs.append(GraphConv(h_feats, num_classes))

    def forward(self, g, h):
        for ell in range(self.num_layers-1):
            h = self.convs[ell](g, h)
            h = F.relu(h)
        h = self.convs[-1](g, h)
        return h
```
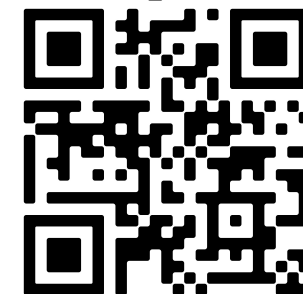


Train F1-score / Iters

Legend:
- Train: GCN (layer=2)
- Valid: GCN (layer=2)
- Train: GCN (layer=4)
- Valid: GCN (layer=4)
- Train: GCN (layer=6)
- Valid: GCN (layer=6)
- Train: GCN (layer=8)
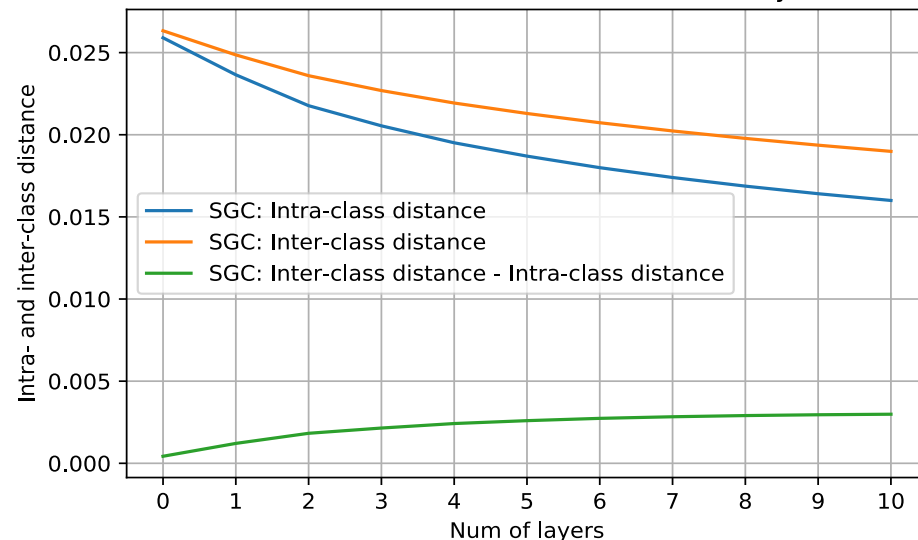- Valid: GCN (layer=8)

# Motivation

- In this paper, we aim at answering two fundamental questions:
    - *Q1: Does increasing depth really impair the expressive power of GCNs?*
    - *Q2: If GCN is expressive, then why do deep GCNs generalize poorly?*

# Q1: *Does increasing depth really impair the expressive power of GCNs?*

- Over-smoothing [1] : a phenomenon where all node embeddings converge to a single vector after applying multiple graph convolution operations to the node features
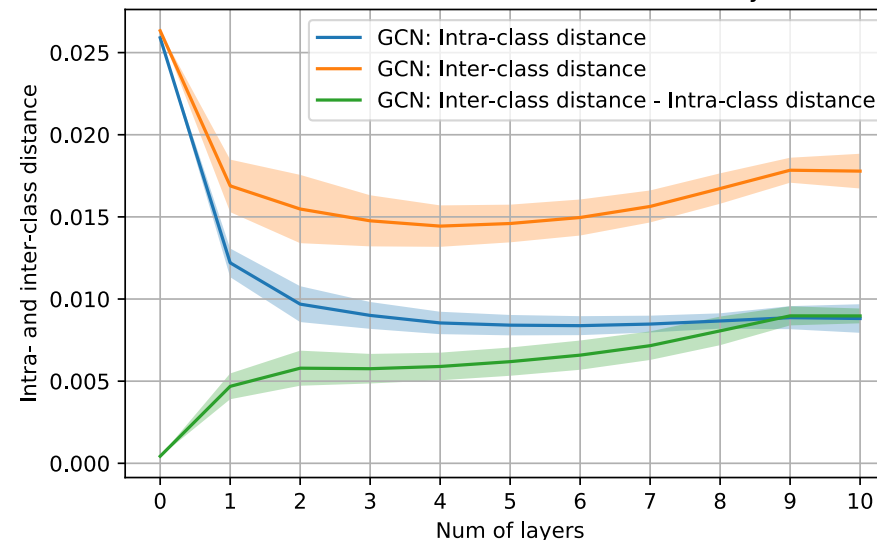
$$\mathbf{H}^{(\ell)} = \mathbf{L}\mathbf{H}^{(\ell-1)}, \mathbf{H}^{(0)} = \mathbf{X}$$

Intra- and inter-class distance / Num of layers

$$\mathbf{H}^{(\ell)} = \sigma(\mathbf{L}\mathbf{H}^{(\ell-1)}\mathbf{W}^{(\ell)}), \mathbf{H}^{(0)} = \mathbf{X}$$

Intra- and inter-class distance / Num of layers

SGC: Intra-class distance
SGC: Inter-class distance
SGC: Inter-class distance - Intra-class distance

GCN: Intra-class distance
GCN: Inter-class distance
GCN: Inter-class distance - Intra-class distance

[1] Li, Qimai, Zhichao Han, and Xiao-Ming Wu. "Deeper insights into graph convolutional networks for semi-supervised learning." Thirty-Second AAAI conference on artificial intelligence. 2018.

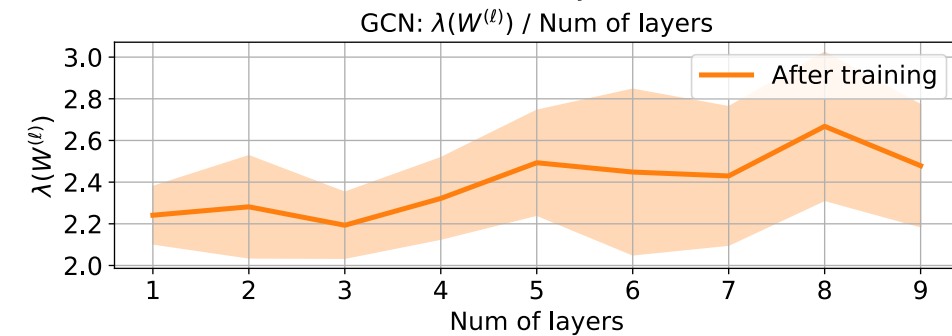# Q1: *Does increasing depth really impair the expressive power of GCNs?*
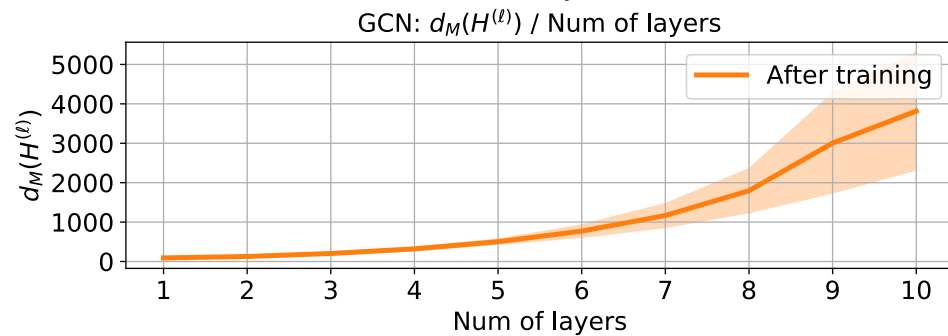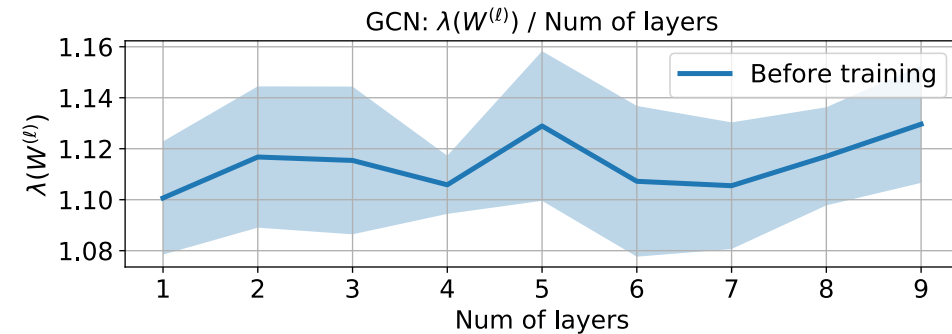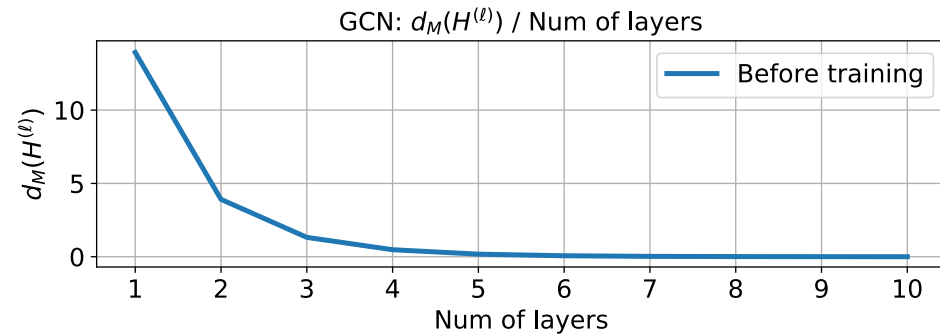
- [2] takes non-linearity and weight matrices into consideration.
- Notations:
  - Expressive power $d_{\mathcal{M}}(\mathbf{H}^{(\ell)})$ as the distance of node embeddings $\mathbf{H}^{(\ell)}$ to a subspace $\mathcal{M}$ that only has node degree information.
  - $\lambda_L$ as the second largest eigenvalue of Laplacian, $\lambda_W$ as the largest singular value of weight matrices
- They show $d_{\mathcal{M}}(\mathbf{H}^{(\ell)}) \leq (\lambda_L \lambda_W)^{\ell} d_{\mathcal{M}}(\mathbf{H}^{(0)})$, i.e., the expressive power will be exponentially decreasing (if $\lambda_L \lambda_W < 1$) or increasing (if $\lambda_L \lambda_W > 1$) as the number of layers increases.

[2] Oono, Kenta, and Taiji Suzuki. "Graph Neural Networks Exponentially Lose Expressive Power for Node Classification." *International Conference on Learning Representations*. 2019.

# Q1: *Does increasing depth really impair the expressive power of GCNs?*

- However, the above assumption (*i.e.*, $\lambda_L \lambda_W < 1$) not always hold.

- For example,
  - Let assume weight matrices $W^{(\ell)} \in \mathbb{R}^{d_{\ell-1} \times d_\ell}$ is initialized by uniform distribution $\mathcal{N}(0, \sqrt{1/d_{\ell-1}})$.
  - By the Gordon's theorem for Gaussian matrices, we know that the expected largest singular value is bounded by $\mathbb{E}[\lambda_W] \leq 1 + \sqrt{d_\ell/d_{\ell-1}}$.
  - This also hold for other initializations.

- Besides, since real-world graphs are sparse, $\lambda_L$ is close to 1.
  - Cora $\lambda_L$=0.9964, Citeseer $\lambda_L$=0.9987, PubMed $\lambda_L$=0.9905

[2] Oono, Kenta, and Taiji Suzuki. "Graph Neural Networks Exponentially Lose Expressive Power for Node Classification." *International Conference on Learning Representations*. 2019.

# Q1: *Does increasing depth really impair the expressive power of GCNs?*

- Besides, we empirically test on real-world dataset

# Q1: *Does increasing depth really impair the expressive power of GCNs?*

- Deeper GCNs have stronger expressive power than the shallow GCNs.
  - [3] shows an appropriately trained GCNs is as expressive as 1-ML test
  - An *L*-layer GCN can encode any different computation tree into different representations.
  - Then, we can characterize the expressiveness of *L*-layer GCN by the number of computation graphs it can encode

**Theorem 1.** *Suppose $\mathcal{T}^L$ is a computation tree with binary node features and node degree at least d. Then the richness of the output of L-GCN defined on $\mathcal{T}^L$ is at least $|L\text{-}GCN(\mathcal{T}^L)| \geq 2(d-1)^{L-1}$.*

[3] Morris, Christopher, et al. "Weisfeiler and leman go neural: Higher-order graph neural networks." *Proceedings of the AAAI Conference on Artificial Intelligence*. Vol. 33. No. 01. 2019.

# Q1: *Does increasing depth really impair the expressive power of GCNs?*

• Besides, we provide global convergence of GCNs

**Theorem 2.** *Let $\boldsymbol{\theta}_t = \{\mathbf{W}_t^{(\ell)} \in \mathbb{R}^{d_{\ell-1} \times d_\ell}\}_{\ell=1}^{L+1}$ be the model parameter at the $t$-th iteration and using square loss $\mathcal{L}(\boldsymbol{\theta}) = \frac{1}{2}\|\mathbf{H}^{(L)}\mathbf{W}^{(L+1)} - \mathbf{Y}\|_{\mathrm{F}}^2$, $\mathbf{H}^{(\ell)} = \sigma(\mathbf{L}\mathbf{H}^{(\ell-1)}\mathbf{W}^{(\ell)})$ as objective function. Then, under the condition that $d_L \geq N$ we can obtain $\mathcal{L}(\boldsymbol{\theta}_T) \leq \epsilon$ if $T \geq C(L)\log(\mathcal{L}(\boldsymbol{\theta}_0)/\epsilon)$, where $\epsilon$ is the desired error and $C(L)$ is a function of GCN depth $L$ that grows as GCN becomes deeper.*

• It is still unclear why a deeper GCN has worse performance than a shallow GCN during the evaluation phase.
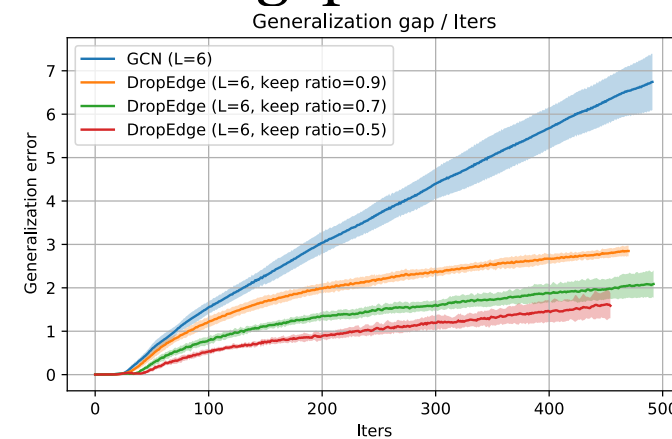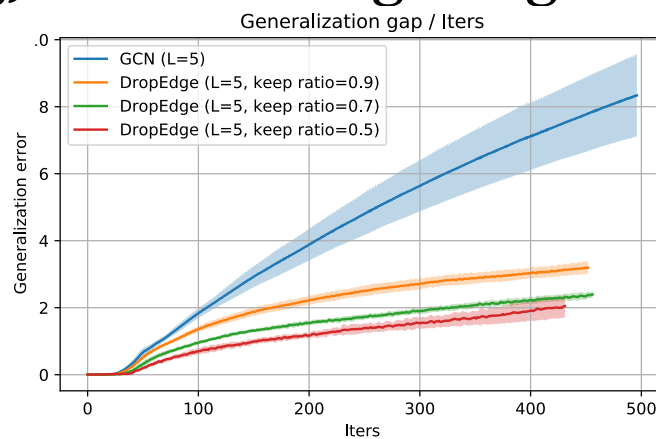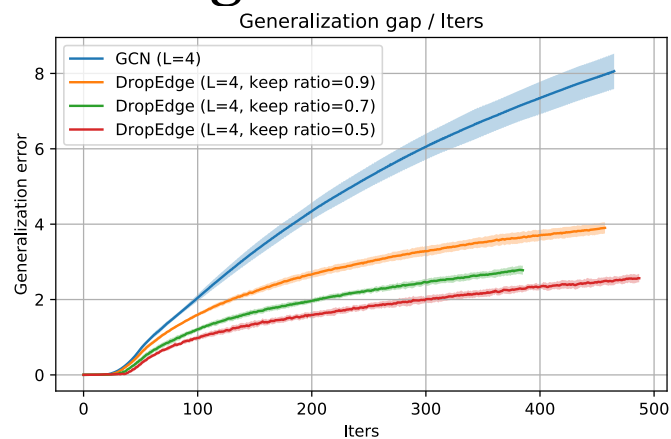
# *Q2:* *If GCN is expressive, why then deep GCNs generalize poorly?*

- To answer this question, we provide a different view by analyzing the impact of GCN structures on the generalization.

- We study the generalization ability of GCNs via *transductive uniform stability*:
  - difference between the training and testing errors for the random partition of a full dataset into training and testing sets.

- Interesting observation:
  - Existing methods that originally designed to alleviate the over-smoothing issue (*e.g., SGC, APPNP, GCNII, DropEdge, PairNorm*) all enjoys a better generalization power than classical GCN.
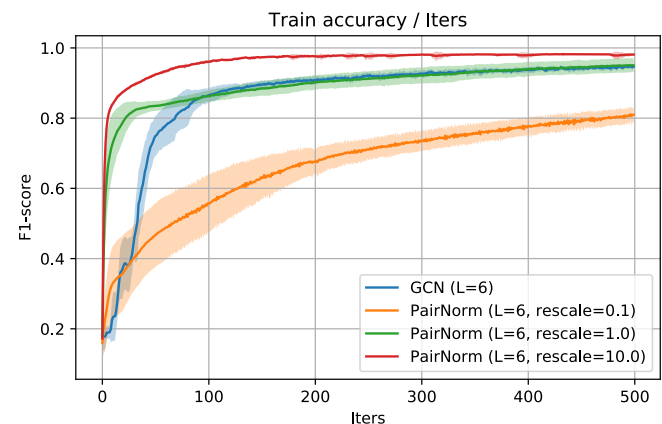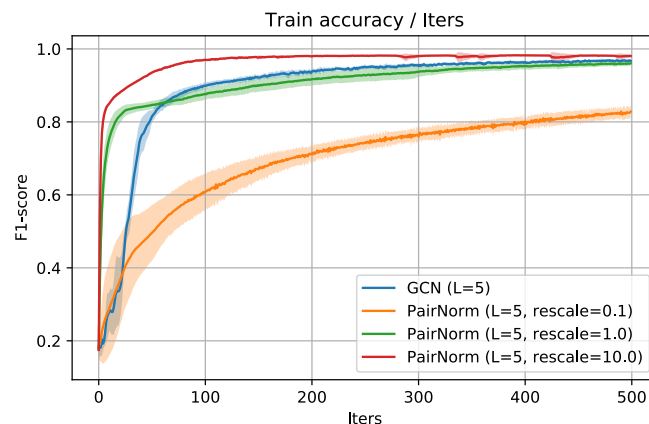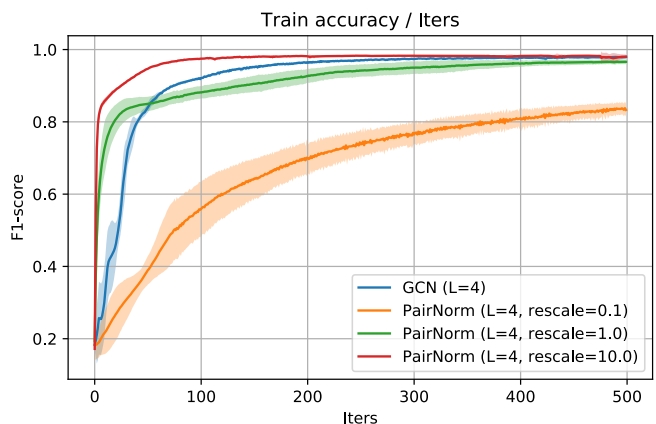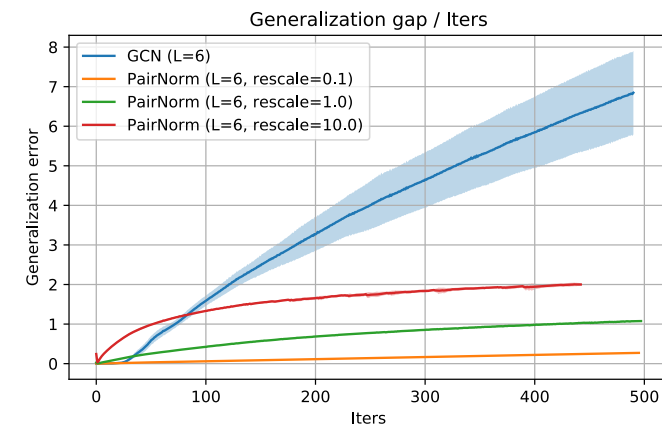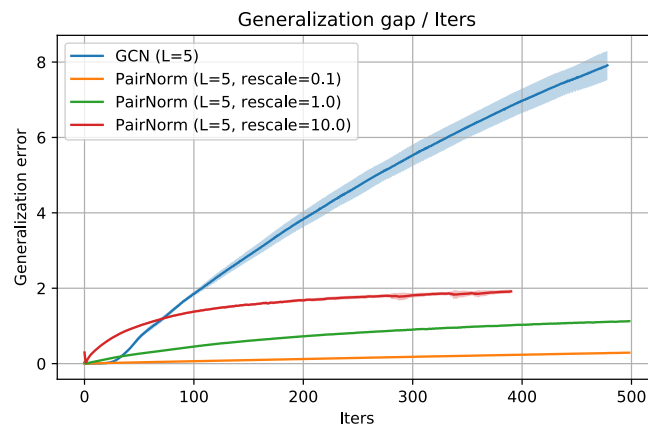
# Q2: *If GCN is expressive, then Why do deep GCNs generalize poorly?*

- For example, **DropEdge** is hurting the training accuracy (i.e., not alleviating over-smoothing) but reducing the generalization gap
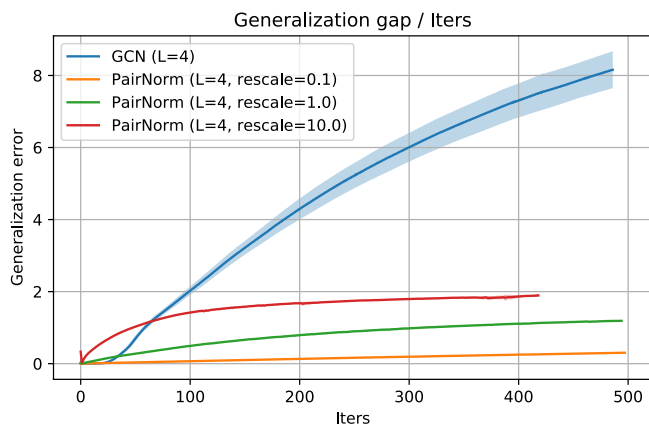
# Q2: *If GCN is expressive, then why do deep GCNs generalize poorly?*

- For example, **PairNorm** is hurting the training accuracy (i.e., not alleviating over-smoothing) but reducing the generalization gap

# *Q2:* *If GCN is expressive, then why do deep GCNs generalize poorly?*

- Informal statement on generalization result

**Theorem 4** (Informal). *We say model is $\epsilon$-uniformly stable with $\epsilon = \frac{2\eta\rho_f G_f}{m}\sum_{t=1}^{T}(1+\eta L_f)^{t-1}$ where the result of $\rho_f, G_f, L_f$ are summarized in Table 1, and other related constants as*

$$B_d^\alpha = (1-\alpha)\sum_{\ell=1}^{L}(\alpha\sqrt{d})^{\ell-1} + (\alpha\sqrt{d})^L, \ B_w^\beta = \beta B_w + (1-\beta),$$
$$B_{\ell,d}^{\alpha,\beta} = \max\left\{\beta\big((1-\alpha)L + \alpha\sqrt{d}\big), (1-\alpha)LB_w^\beta + 1\right\}. \tag{1}$$

Table 1: Comparison of uniform stability constant $\epsilon$ of GCN variants, where $\mathcal{O}(\cdot)$ is used to hide constants that shared between all bounds.

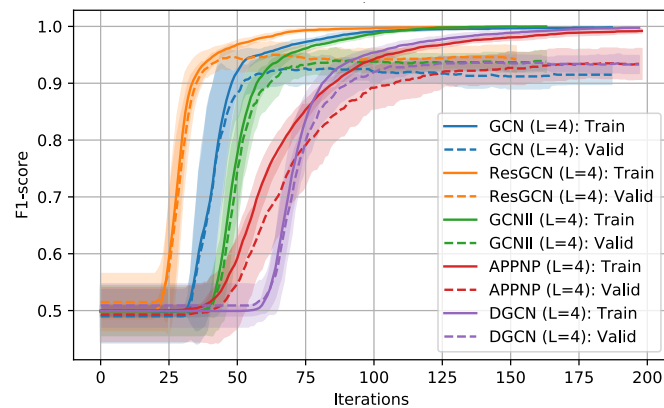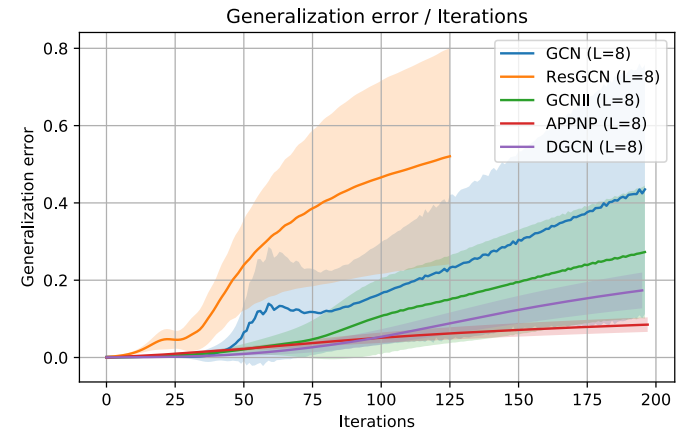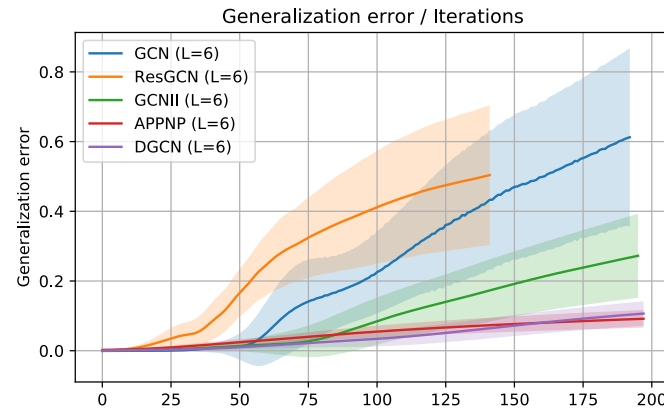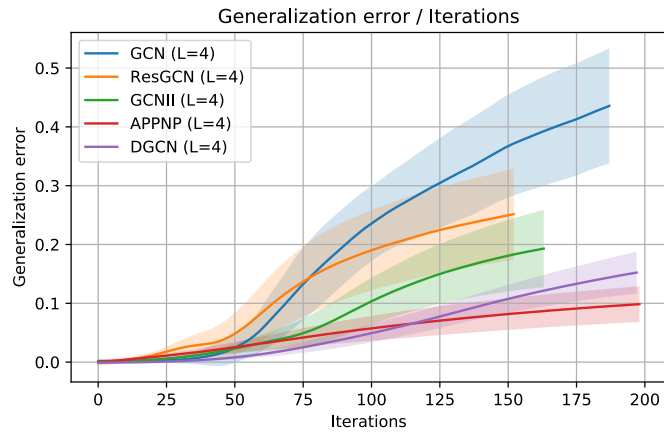| | $\rho_f$ and $G_f$ | $L_f$ | $C_1$ and $C_2$ |
|---|---|---|---|
| $\epsilon_{\text{GCN}}$ | $\mathcal{O}(C_1^L C_2)$ | $\mathcal{O}\big(C_1^L C_2\big((L+2)C_1^L C_2 + 2\big)\big)$ | $C_1 = \max\{1, \sqrt{d}B_w\}, \ C_2 = \sqrt{d}(1+B_x)$ |
| $\epsilon_{\text{ResGCN}}$ | $\mathcal{O}(C_1^L C_2)$ | $\mathcal{O}\big(C_1^L C_2\big((L+2)C_1^L C_2 + 2\big)\big)$ | $C_1 = 1 + \sqrt{d}B_w, \ C_2 = \sqrt{d}(1+B_x)$ |
| $\epsilon_{\text{APPNP}}$ | $\mathcal{O}(C_1)$ | $\mathcal{O}\big(C_1\big(C_1 C_2\big) + 1\big)$ | $C_1 = B_d^\alpha B_x, C_2 = \max\{1, B_w\}$ |
| $\epsilon_{\text{GCNII}}$ | $\mathcal{O}(\beta C_1^L C_2)$ | $\mathcal{O}\big(\alpha\beta C_1^L C_2\big((\alpha\beta L + 2)C_1^L C_2 + 2\beta\big)\big)$ | $C_1 = \max\{1, \alpha\sqrt{d}B_w^\beta\}, \ C_2 = \sqrt{d} + B_{\ell,d}^{\alpha,\beta}B_x$ |
| $\epsilon_{\text{DGCN}}$ | $\mathcal{O}(C_1)$ | $\mathcal{O}(C_1(C_1 C_2) + 1)$ | $C_1 = (\sqrt{d})^L B_x, C_2 = \max\{1, B_w\}$ |

# Proposed GNN architecture

- Based on our generalization analysis, we propose ***Decoupled GCN***, with the following forward propagation rule.
  - $\alpha_\ell, \beta_\ell$ are trainable parameters
  - $\mathbf{P} = \mathbf{D}^{-1/2}\mathbf{A}\mathbf{D}^{-1/2}$ and $\mathbf{P}^\ell$ stands for $\mathbf{P}$ to the power of $\ell$

$$\mathbf{Z} = \sum_{\ell=1}^{L} \alpha_\ell f^{(\ell)}(\mathbf{X}), \ f^{(\ell)}(\mathbf{X}) = \mathbf{P}^\ell \mathbf{X} \left( \beta_\ell \mathbf{W}^{(\ell)} + (1 - \beta_\ell)\mathbf{I} \right))$$

# Empirical validation

- Validate the correctness of the theoretical results on synthetic dataset

Empiri

• Validate t

Train F1-score / Iters

Train F1-score / Iters

Table 3: Comparison of F1-score on OGB-Arxiv dataset for different number of layers

| Model | $\alpha$ | 2 Layers | 4 Layers | 8 Layers | 12 Layers | 16 Layers |
|---|---|---|---|---|---|---|
| **GCN** | $-$ | $71.02\% \pm 0.14$ | $71.56\% \pm 0.19$ | $71.28\% \pm 0.33$ | $70.28\% \pm 0.23$ | $69.37\% \pm 0.46$ |
| **ResGCN** | $-$ | $70.66\% \pm 0.48$ | $72.41\% \pm 0.31$ | $72.56\% \pm 0.31$ | $72.46\% \pm 0.23$ | $72.11\% \pm 0.28$ |
| **GCNII** | $0.9$ | $71.35\% \pm 0.21$ | $72.57\% \pm 0.23$ | $72.06\% \pm 0.42$ | $71.31\% \pm 0.62$ | $69.99\% \pm 0.80$ |
| **GCNII** | $0.8$ | $71.14\% \pm 0.27$ | $72.32\% \pm 0.19$ | $71.90\% \pm 0.41$ | $71.21\% \pm 0.23$ | $70.56\% \pm 0.72$ |
| **GCNII** | $0.5$ | $70.54\% \pm 0.30$ | $72.09\% \pm 0.25$ | $71.92\% \pm 0.32$ | $71.24\% \pm 0.47$ | $71.02\% \pm 0.58$ |
| **APPNP** | $0.9$ | $67.38\% \pm 0.34$ | $68.02\% \pm 0.55$ | $66.62\% \pm 0.48$ | $67.43\% \pm 0.50$ | $67.42\% \pm 1.00$ |
| **APPNP** | $0.8$ | $66.71\% \pm 0.32$ | $68.25\% \pm 0.43$ | $66.40\% \pm 0.89$ | $66.51\% \pm 2.09$ | $66.56\% \pm 0.74$ |
| **DGCN** | $-$ | $71.21\% \pm 0.25$ | $72.29\% \pm 0.18$ | $72.39\% \pm 0.21$ | $\mathbf{72.63\% \pm 0.12}$ | $72.41\% \pm 0.07$ |