

1 We thank all the reviewers for the constructive comments. We’ve done additional experiments to address the concerns  
 2 about the comparison with very recent efforts published this year, and clarified the confusions one-by-one. We reiterate  
 3 our contributions: (1) Global & local metrics; (2) general DGN to improve GNNs even comparing with recent efforts.

4 **To Reviewer 1. Q:** Comparison with GCNII\*<sup>1</sup>. We explore to relieve the over-smoothing issue from the perspective  
 5 of group normalization, and propose a general module – DGN. It can improve different GNN backbones including  
 6 GCNII\*. Our experiments focus on the comparison among existing normalization approaches, rather than fine tuning  
 7 hyperparameters for state-of-the-art performance. The mentioned paper<sup>1</sup> appeared (July 2020) after NeurIPS submission.  
 8 We will include the comparison results on backbone network GCNII\* in the revised version. Table 1 presents part of  
 9 the results, where DGN consistently outperforms the other normalization approaches for different network depths. NN  
 10 (none normalization) denotes the original GCNII\* model. It could be further improved by DGN.

Table 1: The test accuracy in percentage over backbone network GCNII\* on the dataset Cora.

# Layers	NN	BN	PN	DGN	# Layers	NN	BN	PN	DGN	# Layers	NN	BN	PN	DGN
2	80.2	76.5	74.2	<b>83.0</b>	16	83.5	74.4	73.3	<b>85.1</b>	64	85.3	73.7	68.1	<b>86.4</b>

11 **Q:** Metrics. (1) Group distance ratio measures the over-smoothing based on the community impact on the preferred node  
 12 distance, which has not been studied before. This intuitive metric provides us an empirical guidance to design DGN  
 13 module to improve model performance. (2) As mentioned in Section 2.3, both over-fitting and over-smoothing influence  
 14 metric  $G_{Ins}$  of instance information gain. For SGC model, we project node features into label space in the first layer  
 15 and remove all the following trainable matrices. The output representation is generated from the final convolutional  
 16 layer to eliminate the impact of prediction layer. We observe that  $G_{Ins}$  still decreases with more layers, which validates  
 17 the rationality of our metric tailored to measure the mutual information between graph signals. Indeed, GCN/GAT  
 18 usually transforms features into label space in the final convolutional layer, and removes the redundant prediction layer.

19 **Q:** Training process. Thanks for pointing it out. We find that DGN helps GNN models converge to a better training/validation  
 20 loss quickly within hundreds of epochs. Due to space limit, we will include it in the revised version.

21 **Q:** Missing features. The current benchmarks associated informative node attributes usually require shallow GNNs. But  
 22 the over-smoothing issue only appears in the scenarios necessary to stack deeper GNNs to access useful neighbors many  
 23 hops away. As discussed in Section 3.3, one of such scenarios, missing features, exists in many real-world applications  
 24 and has been studied before. Our DGN is shown to effectively tackle this issue and enable deeper neural architectures.  
 25

26 **To Reviewer 2. Q:** Increased params. Thanks for the comment. In GNNs, the increased parameters may instead lead to  
 27 over-fitting. As illustrated in Figure 3 and indicated by group distance ratio, the performance improvement of DGN is  
 28 indeed brought by soft clustering to disentangle unrelated nodes, comparing with the single group in PN & BN.

29 **Q:** Deeper vs shallow models. The deeper model unleashes power in the scenario where a node really needs information  
 30 many hops away to better represent its label. We study such a complex scenario with missing node features in Section  
 31 3.3. As also shown in Table 2 and Appendix, deeper GNNs+DGN achieve better performance than the shallow ones.

32 **Q:** Softmax. We’ve tried Softmax temperatures and regularizations to enforce that each node is assigned to one group.  
 33 We empirically found that soft clustering itself can achieve the similar objective and reduce computation complexity.

34 **To Reviewer 3. Q:** Comparing attentive normalization. Thanks for pointing it out, and we will include the comparison  
 35 in the revised version. Here we present part of results of attentive normalization<sup>2</sup> (AN), attentive context normalization<sup>3</sup>  
 36 (ACN) and our DGN in Table 2, where DGN consistently outperforms them in all cases. This is because ACN has  
 37 only one group and AN additionally samples random groups during model inference. They are not in line with the  
 38 transductive node classification task where the underlying graph has a series of fixed community structures. Furthermore,  
 39 the motivations for our DGN and AN/ACN are different. While AN/ACN target at capturing the long-range relations  
 between pixels, DGN intends to improve the distance between different groups to mitigate the over-smoothing issue.

Table 2: Test accuracy in percentage over GCN model for dataset Cora.

# Layers	AN	ACN	DGN	# Layers	AN	ACN	DGN	# Layers	AN	ACN	DGN
2	77.1	74.4	<b>82.0</b>	15	16.3	42.8	<b>75.2</b>	30	17.2	20.9	<b>73.2</b>

40 **Q:** Zero assignment. A similar case happens in the right part of Figure 3 in our paper, where only 6 out of 10 groups are  
 41 assigned with nodes. In the hyperparameter study, we further find that this issue could be solved by simply choosing a  
 42 proper group number close to the class categories. We will test the regularization term in paper<sup>2</sup> in the revised version.  
 43

44 **To Reviewer 5. Q:** Case studies. We appreciate that you found the work valuable and pointed out this question. As  
 45 mentioned in Section 2.3, SGC achieves better performance with a larger value of  $K$ , since the simplifying mechanism  
 46 avoids over-fitting to enable layer stacking. Compared with Table 1, Table 2 presents results on the complex scenarios  
 47 with missing features. As explained in Section 3.3, the deeper models would be preferred since we need to exploit the  
 48 distant neighbors to collect sparse information. We will explain them in detail in the revised experiment section.

<sup>1</sup>Chen et al. “Simple and Deep Graph Convolutional Networks,” ICML 2020. [Available in arXiv on 4 Jul 2020]

<sup>2</sup>Wang et al. “Attentive Normalization for Conditional Image Generation,” CVPR 2020. [Available in arXiv on 8 Apr 2020]

<sup>3</sup>Sun et al. “ACNe: Attentive Context Normalization for Robust Permutation-Equivariant Learning,” CVPR 2020. [23 Apr 2020]