

1 **To Reviewer #1:** Thanks for your positive comments on our paper!

2 **Q1: Why is Full-Batch outperformed by LADIES?**

3 **A1:** It is true that LADIES is designed as an approximation of original GCN. To answer your question, here we plot  
4 the F1-score of both full-batch GCN and LADIES on the PubMed dataset for 300 epochs without early stop in Figure  
5 1. From Figure 1(a), we can see that, full-batch GCN can achieve higher F1-Score than LADIES on the training set.  
6 Nevertheless, on the validation and test datasets, we can see from Figures 1(b) and 1(c) that LADIES can achieve  
7 significantly higher F1-Score than full-batch GCN. This suggests that LADIES has better generalization performance  
8 than full-batch GCN. The reason is: real graphs are often noisy and incomplete. Full-batch GCN uses the entire graph  
9 in the training phase, which can cause overfitting to the noise. In sharp contrast, LADIES employs stochastic sampling  
10 to use partial information of the graph, and therefore can mitigate the noise of graph and avoid overfitting to the training  
11 data. At a high-level, the sampling scheme adopted in LADIES shares a similar spirit as bagging/bootstrap [1], which is  
12 known to improve the generalization performance of machine learning predictors. We will add these additional plots as  
13 well as the above explanation in the camera-ready version.

14 [1] Leo Breiman. Bagging predictors. Machine learning, 24 (2):123–140, 1996

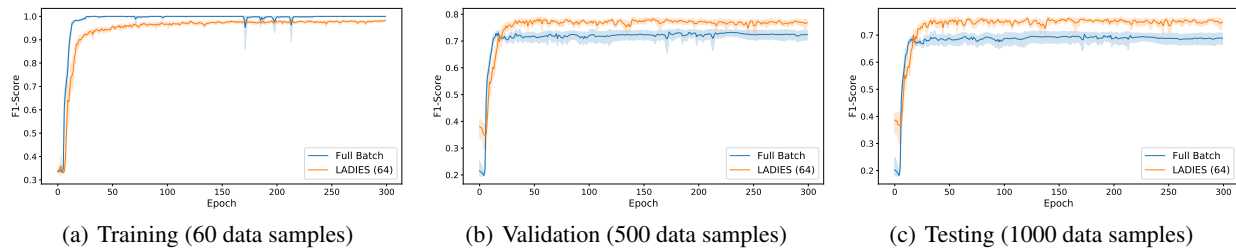


Figure 1: Experiments on the PubMed dataset, which contains 19717 nodes and 44338 edges. We plot the F1-score of both full-batch GCN and LADIES every epoch on (a) Training dataset (b) Validation dataset and (c) Testing dataset.

15 **To Reviewer #2**

16 We would like to highlight the major contributions of our paper as follows: how to efficiently train GCN on large-scale  
17 graphs is a long-standing research problem since the seminal work of GCN (Kipf et al., 2017). Our proposed algorithm  
18 can use a remarkably small sampled node number (64 nodes per GCN layer, which determines the total memory  
19 consumption) to achieve the best known performance with small time and memory consumption, compared with  
20 existing state-of-the-art GCN training algorithms. This suggests that using our algorithm, we can train an accurate GCN  
21 for graphs with very large scale and high density.

22 **Q1: “Similar names generate several misunderstandings and confusions”**

23 **A1:** We apologize the title causes confusing. However, we would like to clarify that our paper entitles “Layer-dependent  
24 importance sampling” rather than “Layered importance sampling”. We would like to emphasize that our “Layer-  
25 dependent importance sampling” is fundamentally different from “Layered importance sampling” discussed in the  
26 paper pointed out by you. We will clarify it in future version of the paper.

27 **Q2: “Importance sampling in the Monte Carlo community has a clear different meaning”**

28 **A2:** We would like to clarify that the terminology of “importance sampling” for GCN training follows the existing work  
29 [2, 3]. Moreover, the “importance sampling” used in our paper belongs to the general category of importance sampling.  
30 As we stated in the paper, we sample a subset of nodes to approximate the exact computation of embeddings in each  
31 layer. In order to reduce the estimation variance, we assign important weights to all candidate nodes (as computed in  
32 (8)) and then apply non-uniformly sampling according to the weights to estimate the embeddings.

33 **To Reviewer #3**

34 **Q1: “The proposed algorithm cannot be considered as a novel improvement in the level of NeurIPS ...”**

35 **A1:** We respectfully disagree with your comment that our algorithm is just a trick in the implementation of FastGCN.  
36 Compared with FastGCN, our algorithm is different in the following aspects: (1) our algorithm restricts the candidate  
37 nodes in the union of the neighborhoods of the sampled nodes in the upper layer, which can lead to a much denser  
38 sampled adjacency matrix compared with FastGCN; (2) our sampling method is layer-dependent, and the importance  
39 sampling probabilities are novelly derived from optimizing on the theoretical derivation of the modified Laplacian  
40 matrix (see (8)), which leads to smaller estimation variance than FastGCN, as shown in Table 2; and (3) our algorithm  
41 uses normalization to the modified Laplacian matrix in each layer, which further stabilizes the forward process and  
42 avoids exploding/vanishing gradient. Because of the above innovative and principled algorithmic designs, our algorithm  
43 enjoys lower time and memory complexities, as well as better predictive performance than FastGCN in both theory and  
44 practice. We believe the improvement of the proposed algorithm is novel and meets the standard of NeurIPS.