

Cluster-GCN: An Efficient Algorithm for Training Deep and Large Graph Convolutional Networks

Wei-Lin Chiang*
National Taiwan University
r06922166@csie.ntu.edu.tw

Xuanqing Liu*
University of California, Los Angeles
xqliu@cs.ucla.edu

Si Si
Google Research
sisidaisy@google.com

Yang Li
Google Research
liyang@google.com

Samy Bengio
Google Research
bengio@google.com

Cho-Jui Hsieh
University of California, Los Angeles
chohsieh@cs.ucla.edu

ABSTRACT

Graph convolutional network (GCN) has been successfully applied to many graph-based applications; however, training a large-scale GCN remains challenging. Current SGD-based algorithms suffer from either a high computational cost that exponentially grows with number of GCN layers, or a large space requirement for keeping the entire graph and the embedding of each node in memory. In this paper, we propose Cluster-GCN, a novel GCN algorithm that is suitable for SGD-based training by exploiting the graph clustering structure. Cluster-GCN works as the following: at each step, it samples a block of nodes that associate with a dense subgraph identified by a graph clustering algorithm, and restricts the neighborhood search within this subgraph. This simple but effective strategy leads to significantly improved memory and computational efficiency while being able to achieve comparable test accuracy with previous algorithms. To test the scalability of our algorithm, we create a new Amazon2M data with 2 million nodes and 61 million edges which is more than 5 times larger than the previous largest publicly available dataset (Reddit). For training a 3-layer GCN on this data, Cluster-GCN is faster than the previous state-of-the-art VR-GCN (1523 seconds vs 1961 seconds) and using much less memory (2.2GB vs 11.2GB). Furthermore, for training 4 layer GCN on this data, our algorithm can finish in around 36 minutes while all the existing GCN training algorithms fail to train due to the out-of-memory issue. Furthermore, Cluster-GCN allows us to train much deeper GCN without much time and memory overhead, which leads to improved prediction accuracy—using a 5-layer Cluster-GCN, we achieve state-of-the-art test F1 score 99.36 on the PPI dataset, while the previous best result was 98.71 by [16]. Our codes are publicly available at https://github.com/google-research/google-research/tree/master/cluster_gcn.

ACM Reference Format:

Wei-Lin Chiang, Xuanqing Liu, Si Si, Yang Li, Samy Bengio, and Cho-Jui Hsieh. 2019. Cluster-GCN: An Efficient Algorithm for Training Deep and

Large Graph Convolutional Networks. In *The 25th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD '19)*, August 4–8, 2019, Anchorage, AK, USA. ACM, New York, NY, USA, 10 pages. <https://doi.org/10.1145/3292500.3330925>

1 INTRODUCTION

Graph convolutional network (GCN) [9] has become increasingly popular in addressing many graph-based applications, including semi-supervised node classification [9], link prediction [17] and recommender systems [15]. Given a graph, GCN uses a graph convolution operation to obtain node embeddings layer by layer—at each layer, the embedding of a node is obtained by gathering the embeddings of its neighbors, followed by one or a few layers of linear transformations and nonlinear activations. The final layer embedding is then used for some end tasks. For instance, in node classification problems, the final layer embedding is passed to a classifier to predict node labels, and thus the parameters of GCN can be trained in an end-to-end manner.

Since the graph convolution operator in GCN needs to propagate embeddings using the interaction between nodes in the graph, this makes training quite challenging. Unlike other neural networks that the training loss can be perfectly decomposed into individual terms on each sample, the loss term in GCN (e.g., classification loss on a single node) depends on a huge number of other nodes, especially when GCN goes deep. Due to the node dependence, GCN’s training is very slow and requires lots of memory – back-propagation needs to store all the embeddings in the computation graph in GPU memory.

Previous GCN Training Algorithms: To demonstrate the need of developing a scalable GCN training algorithm, we first discuss the pros and cons of existing approaches, in terms of 1) memory requirement¹, 2) time per epoch² and 3) convergence speed (loss reduction) per epoch. These three factors are crucial for evaluating a training algorithm. Note that memory requirement directly restricts the scalability of algorithm, and the later two factors combined together will determine the training speed. In the following discussion we denote N to be the number of nodes in the graph, F the embedding dimension, and L the number of layers to analyze classic GCN training algorithms.

- **Full-batch** gradient descent is proposed in the first GCN paper [9]. To compute the full gradient, it requires storing all the

*This work was done during the first and the second author’s internship at Google Research.

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).
KDD '19, August 4–8, 2019, Anchorage, AK, USA
© 2019 Copyright held by the owner/author(s).
ACM ISBN 978-1-4503-6201-6/19/08.
<https://doi.org/10.1145/3292500.3330925>

¹Here we consider the memory for storing node embeddings, which is dense and usually dominates the overall memory usage for deep GCN.

²An epoch means a complete data pass.

intermediate embeddings, leading to $O(NFL)$ memory requirement, which is not scalable. Furthermore, although the time per epoch is efficient, the convergence of gradient descent is slow since the parameters are updated only once per epoch.

[memory: bad; time per epoch: good; convergence: bad]

- **Mini-batch SGD** is proposed in [5]. Since each update is only based on a mini-batch gradient, it can reduce the memory requirement and conduct many updates per epoch, leading to a faster convergence. However, mini-batch SGD introduces a significant computational overhead due to the **neighborhood expansion problem**—to compute the loss on a single node at layer L , it requires that node’s neighbor nodes’ embeddings at layer $L - 1$, which again requires their neighbors’ embeddings at layer $L - 2$ and recursive ones in the downstream layers. This leads to time complexity exponential to the GCN depth. GraphSAGE [5] proposed to use a fixed size of neighborhood samples during back-propagation through layers and FastGCN [1] proposed importance sampling, but the overhead of these methods is still large and will become worse when GCN goes deep.
[memory: good; time per epoch: bad; convergence: good]
- **VR-GCN** [2] proposes to use a **variance reduction technique** to reduce the size of neighborhood sampling nodes. Despite successfully reducing the size of samplings (in our experiments VR-GCN with only 2 samples per node works quite well), it requires storing all the intermediate embeddings of all the nodes in memory, leading to $O(NFL)$ memory requirement. If the number of nodes in the graph increases to millions, the memory requirement for VR-GCN may be too high to fit into GPU.
[memory: bad; time per epoch: good; convergence: good.]

In this paper, we propose a novel GCN training algorithm by exploiting the graph clustering structure. We find that the **efficiency of a mini-batch algorithm** can be characterized by the notion of **“embedding utilization”**, which is proportional **to the number of links between nodes in one batch or within-batch links**. This finding motivates us to design the batches using graph clustering algorithms that aims to construct partitions of nodes so that there are more graph links between nodes in the same partition than nodes in different partitions. Based on the graph clustering idea, we proposed Cluster-GCN, an algorithm to design the batches based on efficient graph clustering algorithms (e.g., METIS [8]). We take this idea further by proposing **a stochastic multi-clustering framework** to improve the convergence of Cluster-GCN. Our strategy leads to huge memory and computational benefits. In terms of memory, we only need to store the node embeddings within the current batch, which is $O(bFL)$ with the batch size b . This is significantly better than VR-GCN and full gradient decent, and slightly better than other SGD-based approaches. In terms of computational complexity, our algorithm achieves the same time cost per epoch with gradient descent and is much faster than neighborhood searching approaches. In terms of the convergence speed, our algorithm is competitive with other SGD-based approaches. Finally, our algorithm is simple to implement since we only compute matrix multiplication and no neighborhood sampling is needed. Therefore for Cluster-GCN, we have [memory: good; time per epoch: good; convergence: good].

We conducted comprehensive experiments on several large-scale graph datasets and made the following contributions:

- Cluster-GCN **achieves the best memory usage** on large-scale graphs, especially on deep GCN. For example, Cluster-GCN uses 5x less memory than VRGCN in a 3-layer GCN model on Amazon2M. Amazon2M is a new graph dataset that we construct to demonstrate the scalability of the GCN algorithms. This dataset contains a amazon product co-purchase graph with more than 2 millions nodes and 61 millions edges.
- Cluster-GCN achieves **a similar training speed with VR-GCN** for shallow networks (e.g., 2 layers) but can be faster than VR-GCN when the network goes deeper (e.g., 4 layers), since our complexity is linear to the number of layers L while VR-GCN’s complexity is exponential to L .
- Cluster-GCN is able to **train a very deep network that has a large embedding size**. Although several previous works show that deep GCN does not give better performance, we found that with proper optimization, deeper GCN could help the accuracy. For example, with a 5-layer GCN, we obtain a new benchmark accuracy 99.36 for PPI dataset, comparing with the highest reported one 98.71 by [16].

Implementation of our proposed method is publicly available.³

2 BACKGROUND

Suppose we are given a graph $G = (\mathcal{V}, \mathcal{E}, A)$, which consists of $N = |\mathcal{V}|$ vertices and $|\mathcal{E}|$ edges such that an edge between any two vertices i and j represents their similarity. The corresponding adjacency matrix A is an $N \times N$ sparse matrix with (i, j) entry equaling to 1 if there is an edge between i and j and 0 otherwise. Also, each node is associated with an F -dimensional feature vector and $X \in \mathbb{R}^{N \times F}$ denotes the feature matrix for all N nodes. An L -layer GCN [9] consists of L graph convolution layers and each of them constructs embeddings for each node by **mixing the embeddings of the node’s neighbors in the graph from the previous layer**:

$$Z^{(l+1)} = A'X^{(l)}W^{(l)}, \quad X^{(l+1)} = \sigma(Z^{(l+1)}), \quad (1)$$

where $X^{(l)} \in \mathbb{R}^{N \times F_l}$ is the embedding at the l -th layer for all the N nodes and $X^{(0)} = X$; A' is the normalized and regularized adjacency matrix and $W^{(l)} \in \mathbb{R}^{F_l \times F_{l+1}}$ is the feature transformation matrix which will be learnt for the downstream tasks. Note that for simplicity we assume the feature dimensions are the same for all layers ($F_1 = \dots = F_L = F$). The activation function $\sigma(\cdot)$ is usually set to be the element-wise ReLU.

Semi-supervised node classification is a popular application of GCN. When using GCN for this application, the goal is to learn weight matrices in (1) by minimizing the loss function:

$$\mathcal{L} = \frac{1}{|\mathcal{Y}_L|} \sum_{i \in \mathcal{Y}_L} \text{loss}(y_i, z_i^L), \quad (2)$$

where \mathcal{Y}_L contains all the labels for the labeled nodes; z_i^L is the i -th row of $Z^{(L)}$ with the ground-truth label to be y_i , indicating the final layer prediction of node i . In practice, a cross-entropy loss is commonly used for node classification in multi-class or multi-label problems.

³https://github.com/google-research/google-research/tree/master/cluster_gcn

Table 1: Time and space complexity of GCN training algorithms. L is number of layers, N is number of nodes, $\|A\|_0$ is number of nonzeros in the adjacency matrix, and F is number of features. For simplicity we assume number of features is fixed for all layers. For SGD-based approaches, b is the batch size and r is the number of sampled neighbors per node. Note that due to the variance reduction technique, VR-GCN can work with a smaller r than GraphSAGE and FastGCN. For memory complexity, LF^2 is for storing $\{W^{(l)}\}_{l=1}^L$ and the other term is for storing embeddings. For simplicity we omit the memory for storing the graph (GCN) or sub-graphs (other approaches) since they are fixed and usually not the main bottleneck.

	GCN [9]	Vanilla SGD	GraphSAGE [5]	FastGCN [1]	VR-GCN [2]	Cluster-GCN
Time complexity	$O(L\ A\ _0F + LNF^2)$	$O(d^L NF^2)$	$O(r^L NF^2)$	$O(rLNF^2)$	$O(L\ A\ _0F + LNF^2 + r^L NF^2)$	$O(L\ A\ _0F + LNF^2)$
Memory complexity	$O(LNF + LF^2)$	$O(bd^L F + LF^2)$	$O(br^L F + LF^2)$	$O(brLF + LF^2)$	$O(LNF + LF^2)$	$O(bLF + LF^2)$

3 PROPOSED ALGORITHM

We first discuss the bottleneck of previous training methods to motivate the proposed algorithm.

In the original paper [9], full gradient descent is used for training GCN, but it suffers from high computational and memory cost. In terms of memory, computing the full gradient of (2) by back-propagation requires storing all the embedding matrices $\{Z^{(l)}\}_{l=1}^L$ which needs $O(NFL)$ space. In terms of convergence speed, since the model is only updated once per epoch, the training requires more epochs to converge.

It has been shown that mini-batch SGD can improve the training speed and memory requirement of GCN in some recent works [1, 2, 5]. Instead of computing the full gradient, SGD only needs to calculate the gradient based on a mini-batch for each update. In this paper, we use $\mathcal{B} \subseteq [N]$ with size $b = |\mathcal{B}|$ to denote a batch of node indices, and each SGD step will compute the gradient estimation

$$\frac{1}{|\mathcal{B}|} \sum_{i \in \mathcal{B}} \nabla \text{loss}(y_i, z_i^{(L)}) \quad (3)$$

to perform an update. Despite faster convergence in terms of epochs, SGD will introduce another computational overhead on GCN training (as explained in the following), which makes it having much slower per-epoch time compared with full gradient descent.

Why does vanilla mini-batch SGD have slow per-epoch time? We consider the computation of the gradient associated with one node i : $\nabla \text{loss}(y_i, z_i^{(L)})$. Clearly, this requires the embedding of node i , which depends on its neighbors' embeddings in the previous layer. To fetch each node i 's neighbor nodes' embeddings, we need to further aggregate each neighbor node's neighbor nodes' embeddings as well. Suppose a GCN has $L + 1$ layers and each node has an average degree of d , to get the gradient for node i , we need to aggregate features from $O(d^L)$ nodes in the graph for one node. That is, we need to fetch information for a node's hop- k ($k = 1, \dots, L$) neighbors in the graph to perform one update. Computing each embedding requires $O(F^2)$ time due to the multiplication with $W^{(l)}$, so in average computing the gradient associated with one node requires $O(d^L F^2)$ time.

Embedding utilization can reflect computational efficiency. If a batch has more than one node, the time complexity is less straightforward since different nodes can have overlapped hop- k neighbors, and the number of embedding computation can be less than the worst case $O(bd^L)$. To reflect the computational efficiency of mini-batch SGD, we define the concept of "embedding utilization" to characterize the computational efficiency. During

the algorithm, if the node i 's embedding at l -th layer $z_i^{(l)}$ is computed and is reused u times for the embedding computations at layer $l + 1$, then we say the embedding utilization of $z_i^{(l)}$ is u . For mini-batch SGD with random sampling, u is very small since the graph is usually large and sparse. Assume u is a small constant (almost no overlaps between hop- k neighbors), then mini-batch SGD needs to compute $O(bd^L)$ embeddings per batch, which leads to $O(bd^L F^2)$ time per update and $O(Nd^L F^2)$ time per epoch.

We illustrate the neighborhood expansion problem in the left panel of Fig. 1. In contrary, full-batch gradient descent has the maximal embedding utilization—each embedding will be reused d (average degree) times in the upper layer. As a consequence, the original full gradient descent [9] only needs to compute $O(NL)$ embeddings per epoch, which means on average only $O(L)$ embedding computation is needed to acquire the gradient of one node.

To make mini-batch SGD work, previous approaches try to restrict the neighborhood expansion size, which however do not improve embedding utilization. GraphSAGE [5] uniformly samples a fixed-size set of neighbors, instead of using a full-neighborhood set. We denote the sample size as r . This leads to $O(r^L)$ embedding computations for each loss term but also makes gradient estimation less accurate. FastGCN [1] proposed an important sampling strategy to improve the gradient estimation. VR-GCN [2] proposed a strategy to store the previous computed embeddings for all the N nodes and L layers and reuse them for unsampled neighbors. Despite the high memory usage for storing all the NL embeddings, we find their strategy very useful and in practice, even for a small r (e.g., 2) can lead to good convergence.

We summarize the time and space complexity in Table 1. Clearly, all the SGD-based algorithms suffer from exponential complexity with respect to the number of layers, and for VR-GCN, even though r can be small, they incur huge space complexity that could go beyond a GPU's memory capacity. In the following, we introduce our Cluster-GCN algorithm, which achieves the best of two worlds—the same time complexity per epoch with full gradient descent and the same memory complexity with vanilla SGD.

3.1 Vanilla Cluster-GCN

Our Cluster-GCN technique is motivated by the following question: In mini-batch SGD updates, can we design a batch and the corresponding computation subgraph to maximize the embedding utilization? We answer this affirmative by connecting the concept of embedding utilization to a clustering objective.

Consider the case that in each batch we compute the embeddings for a set of nodes \mathcal{B} from layer 1 to L . Since the same subgraph $A_{\mathcal{B}, \mathcal{B}}$ (links within \mathcal{B}) is used for each layer of computation, we can

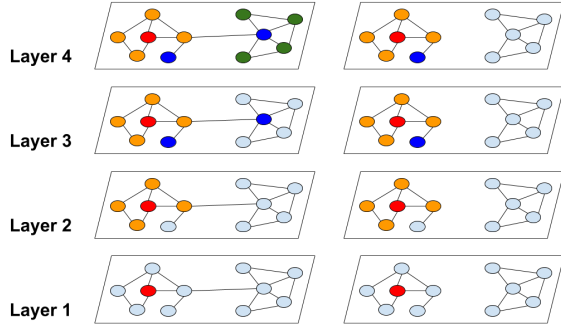


Figure 1: The neighborhood expansion difference between traditional graph convolution and our proposed cluster approach. The red node is the starting node for neighborhood nodes expansion. Traditional graph convolution suffers from exponential neighborhood expansion, while our method can avoid expensive neighborhood expansion.

then see that embedding utilization is the number of edges within this batch $\|A_{\mathcal{B}, \mathcal{B}}\|_0$. Therefore, to maximize embedding utilization, we should design a batch \mathcal{B} to maximize the within-batch edges, by which we connect the efficiency of SGD updates with graph clustering algorithms.

Now we formally introduce **Cluster-GCN**. For a graph G , we partition its nodes into c groups: $\mathcal{V} = [\mathcal{V}_1, \dots, \mathcal{V}_c]$ where \mathcal{V}_t consists of the nodes in the t -th partition. Thus we have c subgraphs as

$$\tilde{G} = [G_1, \dots, G_c] = [\{\mathcal{V}_1, \mathcal{E}_1\}, \dots, \{\mathcal{V}_c, \mathcal{E}_c\}],$$

where each \mathcal{E}_t only consists of the links between nodes in \mathcal{V}_t . After reorganizing nodes, the adjacency matrix is partitioned into c^2 submatrices as

$$A = \tilde{A} + \Delta = \begin{bmatrix} A_{11} & \cdots & A_{1c} \\ \vdots & \ddots & \vdots \\ A_{c1} & \cdots & A_{cc} \end{bmatrix} \quad (4)$$

and

$$\tilde{A} = \begin{bmatrix} A_{11} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & A_{cc} \end{bmatrix}, \Delta = \begin{bmatrix} 0 & \cdots & A_{1c} \\ \vdots & \ddots & \vdots \\ A_{c1} & \cdots & 0 \end{bmatrix}, \quad (5)$$

where each diagonal block A_{tt} is a $|\mathcal{V}_t| \times |\mathcal{V}_t|$ adjacency matrix containing the links within G_t . \tilde{A} is the adjacency matrix for graph \tilde{G} ; A_{st} contains the links between two partitions \mathcal{V}_s and \mathcal{V}_t ; Δ is the matrix consisting of all off-diagonal blocks of A . Similarly, we can partition the feature matrix X and training labels Y according to the partition $[\mathcal{V}_1, \dots, \mathcal{V}_c]$ as $[X_1, \dots, X_c]$ and $[Y_1, \dots, Y_c]$ where X_t and Y_t consist of the features and labels for the nodes in \mathcal{V}_t respectively.

The benefit of this block-diagonal approximation \tilde{G} is that the objective function of GCN becomes decomposable into different batches (clusters). Let \tilde{A}' denotes the normalized version of \tilde{A} , the

final **embedding matrix** becomes

$$Z^{(L)} = \tilde{A}' \sigma(\tilde{A}' \sigma(\cdots \sigma(\tilde{A}' X W^{(0)}) W^{(1)}) \cdots) W^{(L-1)} \quad (6)$$

$$= \begin{bmatrix} \tilde{A}'_{11} \sigma(\tilde{A}'_{11} \sigma(\cdots \sigma(\tilde{A}'_{11} X_1 W^{(0)}) W^{(1)}) \cdots) W^{(L-1)} \\ \vdots \\ \tilde{A}'_{cc} \sigma(\tilde{A}'_{cc} \sigma(\cdots \sigma(\tilde{A}'_{cc} X_c W^{(0)}) W^{(1)}) \cdots) W^{(L-1)} \end{bmatrix}$$

due to the block-diagonal form of \tilde{A} (note that \tilde{A}'_{tt} is the corresponding diagonal block of \tilde{A}'). The loss function can also be decomposed into

$$\mathcal{L}_{\tilde{A}'} = \sum_t \frac{|\mathcal{V}_t|}{N} \mathcal{L}_{\tilde{A}'_{tt}} \quad \text{and} \quad \mathcal{L}_{\tilde{A}'_{tt}} = \frac{1}{|\mathcal{V}_t|} \sum_{i \in \mathcal{V}_t} \text{loss}(y_i, z_i^{(L)}). \quad (7)$$

The Cluster-GCN is then based on the decomposition form in (6) and (7). At each step, we sample a cluster \mathcal{V}_t and then conduct SGD to update based on the gradient of $\mathcal{L}_{\tilde{A}'_{tt}}$, and this only requires the sub-graph A_{tt} , the X_t , Y_t on the current batch and the models $\{W^{(l)}\}_{l=1}^L$. The implementation only requires forward and backward propagation of matrix products (one block of (6)) that is much easier to implement than the neighborhood search procedure used in previous SGD-based training methods.

We use graph clustering algorithms to partition the graph. Graph clustering methods such as **Metis** [8] and **Graclus** [4] aim to **construct the partitions over the vertices in the graph such that within-clusters links** are much more than between-cluster links to better capture the clustering and community structure of the graph. These are exactly what we need because: 1) As mentioned before, the embedding utilization is equivalent to the within-cluster links for each batch. Intuitively, each node and its neighbors are usually located in the same cluster, therefore after a few hops, neighborhood nodes with a high chance are still in the same cluster. 2) Since we replace A by its block diagonal approximation \tilde{A} and the error is proportional to between-cluster links Δ , we need to find a partition to minimize number of between-cluster links.

In Figure 1, we illustrate the neighborhood expansion with full graph G and the graph with clustering partition \tilde{G} . We can see that cluster-GCN can avoid heavy neighborhood search and focus on the neighbors within each cluster. In Table 2, we show two different node partition strategies: random partition versus clustering partition. We partition the graph into 10 parts by using random partition and METIS. Then use one partition as a batch to perform a SGD update. We can see that with the same number of epochs, using clustering partition can achieve higher accuracy. This shows using graph clustering is important and partitions should not be formed randomly.

Time and space complexity. Since each node in \mathcal{V}_t only links to nodes inside \mathcal{V}_t , each node does not need to perform neighborhoods searching outside A_{tt} . The computation for each batch will purely be matrix products $\tilde{A}'_{tt} X_t^{(l)} W^{(l)}$ and some element-wise operations, so the overall time complexity per batch is $O(\|A_{tt}\|_0 F + b F^2)$. Thus the overall time complexity per epoch becomes $O(\|A\|_0 F + N F^2)$. In average, each batch only requires **computing $O(bL)$ embeddings**, which is linear instead of exponential to L . In terms of space complexity, in each batch, we only need to load b samples and store their embeddings on each layer, resulting in **$O(bLF)$ memory** for storing embeddings. Therefore our algorithm is also more memory

Table 2: Random partition versus clustering partition of the graph (trained on mini-batch SGD). Clustering partition leads to better performance (in terms of test F1 score) since it removes less between-partition links. These three datasets are all public GCN datasets. We will explain PPI data in the experiment part. Cora has 2,708 nodes and 13,264 edges, and Pubmed has 19,717 nodes and 108,365 edges.

Dataset	random partition	clustering partition
Cora	78.4	82.5
Pubmed	78.9	79.9
PPI	68.1	92.9

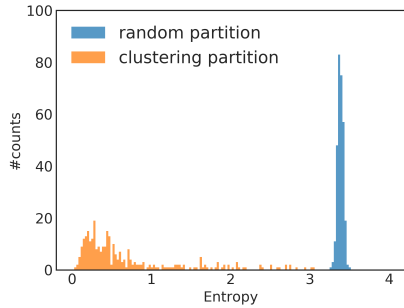


Figure 2: Histograms of entropy values based on the label distribution. Here we present within each batch using random partition versus clustering partition. Most clustering partitioned batches have low label entropy, indicating skewed label distribution within each batch. In comparison, random partition will lead to larger label entropy within a batch although it is less efficient as discussed earlier. We partition the Reddit dataset with 300 clusters in this example.

efficient than all the previous algorithms. Moreover, our algorithm only requires loading a subgraph into GPU memory instead of the full graph (though graph is usually not the memory bottleneck). The detailed time and memory complexity are summarized in Table 1.

3.2 Stochastic Multiple Partitions

Although vanilla Cluster-GCN achieves good computational and memory complexity, there are still two potential issues:

- After the graph is partitioned, **some links** (the Δ part in Eq. (4)) **are removed**. Thus the performance could be affected.
- Graph clustering algorithms tend to **bring similar nodes together**. Hence **the distribution** of a cluster could be **different** from the original data set, leading to a biased estimation of the full gradient while performing SGD updates.

In Figure 2, we demonstrate an example of unbalanced label distribution by using the Reddit data with clusters formed by Metis. We calculate the entropy value of each cluster based on its label distribution. Comparing with random partitioning, we clearly see that **entropy of most clusters are smaller**, indicating that the label **distributions of clusters are biased towards some specific labels**. This increases the variance across different batches and may affect the convergence of SGD.

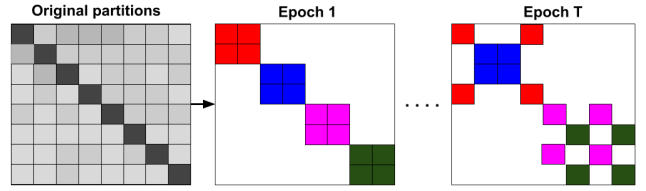


Figure 3: The proposed stochastic multiple partitions scheme. In each epoch, we randomly sample q clusters ($q = 2$ is used in this example) and their between-cluster links to form a new batch. **Same color blocks are in the same batch.**

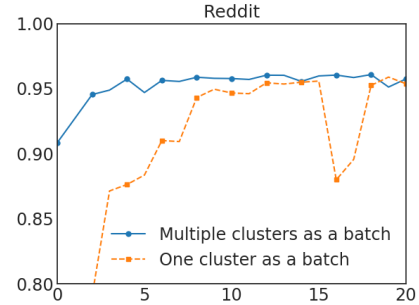


Figure 4: Comparisons of choosing one cluster versus multiple clusters. The former uses 300 partitions. The latter uses 1500 and randomly select 5 to form one batch. We present epoch (x-axis) versus F1 score (y-axis).

To address the above issues, we propose a **stochastic multiple clustering** approach to **incorporate between-cluster links** and reduce **variance across batches**. We first partition the graph into p clusters $\mathcal{V}_1, \dots, \mathcal{V}_p$ with a relatively large p . When constructing a batch B for an SGD update, instead of considering only one cluster, we randomly choose q clusters, denoted as t_1, \dots, t_q and include their nodes $\{\mathcal{V}_{t_1} \cup \dots \cup \mathcal{V}_{t_q}\}$ into the batch. Furthermore, the links between the chosen clusters,

$$\{A_{ij} \mid i, j \in t_1, \dots, t_q\},$$

are added back. In this way, those between-cluster links are re-incorporated and the combinations of clusters make the variance across batches smaller. Figure 3 illustrates our algorithm—for each epochs, different combinations of clusters are chosen as a batch. We conduct an experiment on Reddit to demonstrate the effectiveness of the proposed approach. In Figure 4, we can observe that using multiple clusters as one batch could improve the convergence. Our final Cluster-GCN algorithm is presented in Algorithm 1.

3.3 Issues of training deeper GCNs

Previous attempts of training deeper GCNs [9] seem to suggest that adding more layers is not helpful. However, the datasets used in the experiments may be too small to make a proper justification. For example, [9] considered a graph with only a few hundreds of training nodes for which overfitting can be an issue. Moreover, we observe that the optimization of deep GCN models becomes difficult as it may **impede the information from the first few layers** being passed through. In [9], they adopt a technique similar to **residual**

Algorithm 1: Cluster GCN**Input:** Graph A , feature X , label Y ;**Output:** Node representation \tilde{X}

- 1 Partition graph nodes into c clusters $\mathcal{V}_1, \mathcal{V}_2, \dots, \mathcal{V}_c$ by METIS;
- 2 **for** $iter = 1, \dots, max_iter$ **do**
- 3 Randomly choose q clusters, t_1, \dots, t_q from \mathcal{V} without replacement;
- 4 Form the subgraph \tilde{G} with nodes $\tilde{\mathcal{V}} = [\mathcal{V}_{t_1}, \mathcal{V}_{t_2}, \dots, \mathcal{V}_{t_q}]$ and links $A_{\tilde{\mathcal{V}}, \tilde{\mathcal{V}}}$;
- 5 Compute $g \leftarrow \nabla \mathcal{L}_{A_{\tilde{\mathcal{V}}, \tilde{\mathcal{V}}}}$ (loss on the subgraph $A_{\tilde{\mathcal{V}}, \tilde{\mathcal{V}}}$);
- 6 Conduct Adam update using gradient estimator g
- 7 **Output:** $\{W_l\}_{l=1}^L$

connections [6] to enable the model to carry the information from a previous layer to a next layer. Specifically, they modify (1) to add the hidden representations of layer l into the next layer.

$$X^{(l+1)} = \sigma(A'X^{(l)}W^{(l)}) + X^{(l)} \quad (8)$$

Here we propose another simple technique to improve the training of deep GCNs. In the original GCN settings, each node aggregates the representation of its neighbors from the previous layer. However, under the setting of deep GCNs, the strategy may not be suitable as it does not take the number of layers into account. Intuitively, **neighbors nearby should contribute more** than distant nodes. We thus propose a technique to better address this issue. The idea is to **amplify the diagonal parts of the adjacency matrix A** used in each GCN layer. In this way, we are putting more weights on the representation from the previous layer in the aggregation of each GCN layer. An example is to add an identity to A as follows.

$$X^{(l+1)} = \sigma((A' + I)X^{(l)}W^{(l)}) \quad (9)$$

While (9) seems to be reasonable, using the same weight for all the nodes regardless of their numbers of neighbors may not be suitable. Moreover, it may suffer from numerical instability as values can grow exponentially when more layers are used. Hence we propose a modified version of (9) to better maintain the neighborhoods information and numerical ranges. We first add an identity to the original A and perform the normalization

$$\tilde{A} = (D + I)^{-1}(A + I), \quad (10)$$

and then consider

$$X^{(l+1)} = \sigma((\tilde{A} + \lambda \text{diag}(\tilde{A}))X^{(l)}W^{(l)}). \quad (11)$$

Experimental results of adopting the “diagonal enhancement” techniques are presented in Section 4.3 where we show that this new normalization strategy can help to build deep GCN and achieve SOTA performance.

4 EXPERIMENTS

We evaluate our proposed method for training GCN on two tasks: multi-label and multi-class classification on four public datasets. The statistic of the data sets are shown in Table 3. Note that the Reddit dataset is the largest public dataset we have seen so far for

Table 3: Data statistics

Datasets	Task	#Nodes	#Edges	#Labels	#Features
PPI	multi-label	56,944	818,716	121	50
Reddit	multi-class	232,965	11,606,919	41	602
Amazon	multi-label	334,863	925,872	58	N/A
Amazon2M	multi-class	2,449,029	61,859,140	47	100

Table 4: The parameters used in the experiments.

Datasets	#hidden units	# partitions	#clusters per batch
PPI	512	50	1
Reddit	128	1500	20
Amazon	128	200	1
Amazon2M	400	15000	10

GCN, and the Amazon2M dataset is collected by ourselves and is much larger than Reddit (see more details in Section 4.2).

We include the following state-of-the-art GCN training algorithms in our comparisons:

- Cluster-GCN (Our proposed algorithm): the proposed fast GCN training method.
- VRGCN⁴ [2]: It maintains the historical embedding of all the nodes in the graph and expands to only a few neighbors to speedup training. The number of sampled neighbors is set to be 2 as suggested in [2]⁵.
- GraphSAGE⁶ [5]: It samples a fixed number of neighbors per node. We use the default settings of sampled sizes for each layer ($S_1 = 25, S_2 = 10$) in GraphSAGE.

We implement our method in PyTorch [13]. For the other methods, we use all the original papers’ code from their github pages. Since [9] has difficulty to scale to large graphs, we do not compare with it here. Also as shown in [2] that VRGCN is faster than FastGCN, so we do not compare with FastGCN here. For all the methods we use the Adam optimizer with learning rate as 0.01, dropout rate as 20%, weight decay as zero. The mean aggregator proposed by [5] is adopted and the number of hidden units is the same for all methods. Note that techniques such as (11) is not considered here. In each experiment, we consider the same GCN architecture for all methods. For VRGCN and GraphSAGE, we follow the settings provided by the original papers and set the batch sizes as 512. For Cluster-GCN, the number of partitions and clusters per batch for each dataset are listed in Table 4. Note that clustering is seen as a preprocessing step and its running time is not taken into account in training. In Section 6, we show that graph clustering only takes a small portion of preprocessing time. All the experiments are conducted on a machine with a NVIDIA Tesla V100 GPU (16 GB memory), 20-core Intel Xeon CPU (2.20 GHz), and 192 GB of RAM.

4.1 Training Performance for median size datasets

Training Time vs Accuracy: First we compare our proposed method with other methods in terms of training speed. In Figure 6, the x -axis shows the training time in seconds, and y -axis shows the accuracy (F1 score) on the validation sets. We plot the training time

⁴GitHub link: https://github.com/thu-ml/stochastic_gcnn

⁵Note that we also tried the default sample size 20 in VRGCN package but it performs much worse than sample size= 2.

⁶GitHub link: <https://github.com/williamleif/GraphSAGE>

Table 5: Comparisons of memory usages on different datasets. Numbers in the brackets indicate the size of hidden units used in the model.

	2-layer			3-layer			4-layer		
	VRGCN	Cluster-GCN	GraphSAGE	VRGCN	Cluster-GCN	GraphSAGE	VRGCN	Cluster-GCN	GraphSAGE
PPI (512)	258 MB	39 MB	51 MB	373 MB	46 MB	71 MB	522 MB	55 MB	85 MB
Reddit (128)	259 MB	284 MB	1074 MB	372 MB	285 MB	1075 MB	515 MB	285 MB	1076 MB
Reddit (512)	1031 MB	292 MB	1099 MB	1491 MB	300 MB	1115 MB	2064 MB	308 MB	1131 MB
Amazon (128)	1188 MB	703 MB	N/A	1351 MB	704 MB	N/A	1515 MB	705 MB	N/A

Table 6: Benchmarking on the Sparse Tensor operations in PyTorch and TensorFlow. A network with two linear layers is used and the timing includes forward and backward operations. Numbers in the brackets indicate the size of hidden units in the first layer. Amazon data is used.

	PyTorch	TensorFlow
Avg. time per epoch (128)	8.81s	2.53s
Avg. time per epoch (512)	45.08s	7.13s

versus accuracy for three datasets with 2,3,4 layers of GCN. Since GraphSAGE is slower than VRGCN and our method, the curves for GraphSAGE only appear for PPI and Reddit datasets. We can see that our method is the fastest for both PPI and Reddit datasets for GCNs with different numbers of layers.

For Amazon data, since nodes’ features are not available, an identity matrix is used as the feature matrix X . Under this setting, the shape of parameter matrix $W^{(0)}$ becomes 334863x128. Therefore, the computation is dominated by sparse matrix operations such as $AW^{(0)}$. Our method is still faster than VRGCN for 3-layer case, but slower for 2-layer and 4-layer ones. The reason may come from the speed of sparse matrix operations from different frameworks. VRGCN is implemented in TensorFlow, while Cluster-GCN is implemented in PyTorch whose sparse tensor support are still in its very early stage. In Table 6, we show the time for TensorFlow and PyTorch to do forward/backward operations on Amazon data, and a simple two-layer network are used for benchmarking both frameworks. We can clearly see that TensorFlow is faster than PyTorch. The difference is more significant when the number of hidden units increases. This may explain why Cluster-GCN has longer training time in Amazon dataset.

Memory usage comparison: For training large-scale GCNs, besides training time, memory usage needed for training is often more important and will directly restrict the scalability. The memory usage includes the memory needed for training the GCN for many epochs. As discussed in Section 3, to speedup training, VRGCN needs to save historical embeddings during training, so it needs much more memory for training than Cluster-GCN. GraphSAGE also has higher memory requirement than Cluster-GCN due to the exponential neighborhood growing problem. In Table 5, we compare our memory usage with VRGCN’s memory usage for GCN with different layers. When increasing the number of layers, Cluster-GCN’s memory usage does not increase a lot. The reason is that when increasing one layer, the extra variable introduced is the weight matrix $W^{(L)}$, which is relatively small comparing to the sub-graph and node features. While VRGCN needs to save each layer’s history embeddings, and the embeddings are usually dense and will soon dominate the memory usage. We can see from Table 5

Table 7: The most common categories in Amazon2M.

Categories	number of products
Books	668,950
CDs & Vinyl	172,199
Toys & Games	158,771

that Cluster-GCN is much more memory efficient than VRGCN. For instance, on Reddit data to train a 4-layer GCN with hidden dimension to be 512, VRGCN needs 2064MB memory, while Cluster-GCN only uses 308MB memory.

4.2 Experimental results on Amazon2M

A new GCN dataset: Amazon2M. By far the largest public data for testing GCN is Reddit dataset with the statistics shown in Table 3, which contains about 200K nodes. As shown in Figure 6 GCN training on this data can be finished within a few hundreds seconds. To test the scalability of GCN training algorithms, we constructed a much larger graph with over 2 millions of nodes and 61 million edges based on Amazon co-purchasing networks [11, 12]. The raw co-purchase data is from AMAZON-3M⁷. In the graph, each node is a product, and the graph link represents whether two products are purchased together. Each node feature is generated by extracting bag-of-word features from the product descriptions followed by Principal Component Analysis [7] to reduce the dimension to be 100. In addition, we use the top-level categories as the labels for that product/node (see Table 7 for the most common categories). The detailed statistics of the data set are listed in Table 3.

In Table 8, we compare with VRGCN for GCNs with a different number of layers in terms of training time, memory usage, and test accuracy (F1 score). As can be seen from the table that 1) VRGCN is faster than Cluster-GCN with 2-layer GCN but slower than Cluster-GCN when increasing one layer while achieving similar accuracy. 2) In terms of memory usage, VRGCN is using much more memory than Cluster-GCN (5 times more for 3-layer case), and it is running out of memory when training 4-layer GCN, while Cluster-GCN does not need much additional memory when increasing the number of layers, and achieves the best accuracy for this data when training a 4-layer GCN.

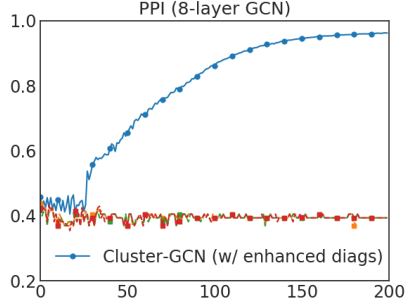
4.3 Training Deeper GCN

In this section we consider GCNs with more layers. We first show the timing comparisons of Cluster-GCN and VRGCN in Table 9. PPI is used for benchmarking and we run 200 epochs for both methods. We observe that the running time of VRGCN grows exponentially

⁷<http://manikvarma.org/downloads/XC/XMLRepository.html>

Table 8: Comparisons of running time, memory and testing accuracy (F1 score) for Amazon2M.

	Time		Memory		Test F1 score	
	VRGCN	Cluster-GCN	VRGCN	Cluster-GCN	VRGCN	Cluster-GCN
Amazon2M (2-layer)	337s	1223s	7476 MB	2228 MB	89.03	89.00
Amazon2M (3-layer)	1961s	1523s	11218 MB	2235 MB	90.21	90.21
Amazon2M (4-layer)	N/A	2289s	OOM	2241 MB	N/A	90.41

**Figure 5: Convergence figure on a 8-layer GCN. We present numbers of epochs (x-axis) versus validation accuracy (y-axis). All methods except for the one using (11) fail to converge.****Table 9: Comparisons of running time when using different numbers of GCN layers. We use PPI and run both methods for 200 epochs.**

	2-layer	3-layer	4-layer	5-layer	6-layer
Cluster-GCN	52.9s	82.5s	109.4s	137.8s	157.3s
VRGCN	103.6s	229.0s	521.2s	1054s	1956s

because of its expensive neighborhood finding, while the running time of Cluster-GCN only grows linearly.

Next we investigate whether using deeper GCNs obtains better accuracy. In Section 4.3, we discuss different strategies of modifying the adjacency matrix A to facilitate the training of deep GCNs. We apply the diagonal enhancement techniques to deep GCNs and run experiments on PPI. Results are shown in Table 11. For the case of 2 to 5 layers, the accuracy of all methods increases with more layers added, suggesting that deeper GCNs may be useful. However, when 7 or 8 GCN layers are used, the first three methods fail to converge within 200 epochs and get a dramatic loss of accuracy. A possible reason is that the optimization for deeper GCNs becomes more difficult. We show a detailed convergence of a 8-layer GCN in Figure 5. With the proposed diagonal enhancement technique (11), the convergence can be improved significantly and similar accuracy can be achieved.

State-of-the-art results by training deeper GCNs. With the design of Cluster-GCN and the proposed normalization approach, we now have the ability for training much deeper GCNs to achieve better accuracy (F1 score). We compare the testing accuracy with other existing methods in Table 10. For PPI, Cluster-GCN can achieve the state-of-art result by training a 5-layer GCN with 2048 hidden units. For Reddit, a 4-layer GCN with 128 hidden units is used.

Table 10: State-of-the-art performance of testing accuracy reported in recent papers.

	PPI	Reddit
FastGCN [1]	N/A	93.7
GraphSAGE [5]	61.2	95.4
VR-GCN [2]	97.8	96.3
GaAN [16]	98.71	96.36
GAT [14]	97.3	N/A
GeniePath [10]	98.5	N/A
Cluster-GCN	99.36	96.60

5 CONCLUSION

We present ClusterGCN, a new GCN training algorithm that is fast and memory efficient. Experimental results show that this method can train very deep GCN on large-scale graph, for instance on a graph with over 2 million nodes, the training time is less than an hour using around 2G memory and achieves accuracy of 90.41 (F1 score). Using the proposed approach, we are able to successfully train much deeper GCNs, which achieve state-of-the-art test F1 score on PPI and Reddit datasets.

Acknowledgement CJH acknowledges the support of NSF via IIS-1719097, Intel faculty award, Google Cloud and Nvidia.

REFERENCES

- [1] Jie Chen, Tengfei Ma, and Cao Xiao. 2018. FastGCN: Fast Learning with Graph Convolutional Networks via Importance Sampling. In *ICLR*.
- [2] Jianfei Chen, Jun Zhu, and Song Le. 2018. Stochastic Training of Graph Convolutional Networks with Variance Reduction. In *ICML*.
- [3] Hanjun Dai, Zornitsa Kozareva, Bo Dai, Alex Smola, and Le Song. 2018. Learning Steady-States of Iterative Algorithms over Graphs. In *ICML*. 1114–1122.
- [4] Inderjit S. Dhillon, Yuqiang Guan, and Brian Kulis. 2007. Weighted Graph Cuts Without Eigenvectors A Multilevel Approach. *IEEE Trans. Pattern Anal. Mach. Intell.* 29, 11 (2007), 1944–1957.
- [5] William L. Hamilton, Rex Ying, and Jure Leskovec. 2017. Inductive Representation Learning on Large Graphs. In *NIPS*.
- [6] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep Residual Learning for Image Recognition. *CVPR* (2016), 770–778.
- [7] H. Hotelling. 1933. Analysis of a complex of statistical variables into principal components. *Journal of Educational Psychology* 24, 6 (1933), 417–441.
- [8] George Karypis and Vipin Kumar. 1998. A fast and high quality multilevel scheme for partitioning irregular graphs. *SIAM J. Sci. Comput.* 20, 1 (1998), 359–392.
- [9] Thomas N. Kipf and Max Welling. 2017. Semi-Supervised Classification with Graph Convolutional Networks. In *ICLR*.
- [10] Ziqi Liu, Chaochao Chen, Longfei Li, Jun Zhou, Xiaolong Li, Le Song, and Yuan Qi. 2019. GeniePath: Graph Neural Networks with Adaptive Receptive Paths. In *AAAI*.
- [11] Julian McAuley, Rahul Pandey, and Jure Leskovec. 2015. Inferring Networks of Substitutable and Complementary Products. In *KDD*.
- [12] Julian McAuley, Christopher Targett, Qinfeng Shi, and Anton van den Hengel. 2015. Image-Based Recommendations on Styles and Substitutes. In *SIGIR*.
- [13] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. 2017. Automatic differentiation in PyTorch. In *NIPS-W*.
- [14] Petar Velićković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. 2018. Graph Attention Networks. (2018).
- [15] Rex Ying, Ruining He, Kaifeng Chen, Pong Eksombatchai, William L. Hamilton, and Jure Leskovec. 2018. Graph Convolutional Neural Networks for Web-Scale Recommender Systems. In *KDD*.

Table 11: Comparisons of using different diagonal enhancement techniques. For all methods, we present the best validation accuracy achieved in 200 epochs. PPI is used and dropout rate is 0.1 in this experiment. Other settings are the same as in Section 4.1. The numbers marked red indicate poor convergence.

	2-layer	3-layer	4-layer	5-layer	6-layer	7-layer	8-layer
Cluster-GCN with (1)	90.3	97.6	98.2	98.3	94.1	65.4	43.1
Cluster-GCN with (10)	90.2	97.7	98.1	98.4	42.4	42.4	42.4
Cluster-GCN with (10) + (9)	84.9	96.0	97.1	97.6	97.3	43.9	43.8
Cluster-GCN with (10) + (11), $\lambda = 1$	89.6	97.5	98.2	98.3	98.0	97.4	96.2

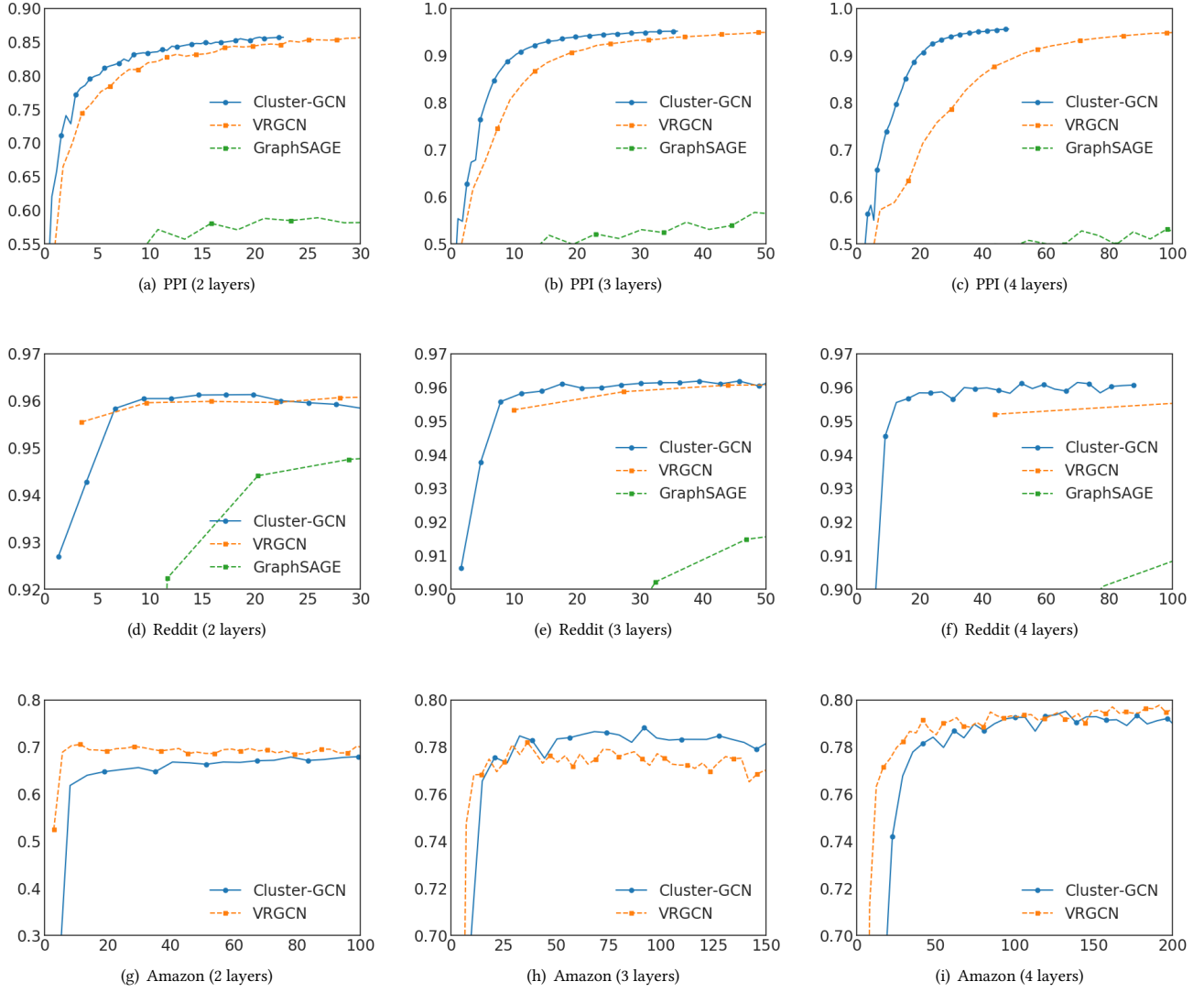


Figure 6: Comparisons of different GCN training methods. We present the relation between training time in seconds (x-axis) and the validation F1 score (y-axis).

- [16] Jiani Zhang, Xingjian Shi, Junyuan Xie, Hao Ma, Irwin King, and Dit-Yan Yeung. 2018. GaAN: Gated Attention Networks for Learning on Large and Spatiotemporal Graphs. In *UAI*.
- [17] Muhan Zhang and Yixin Chen. 2018. Link Prediction Based on Graph Neural Networks. In *NIPS*.

Table 12: The training, validation, and test splits used in the experiments. Note that for the two amazon datasets we only split into training and test sets.

Datasets	Task	Data splits (Tr./Val./Te.)
PPI	Inductive	44906/6514/5524
Reddit	Inductive	153932/23699/55334
Amazon	Inductive	91973/242890
Amazon2M	Inductive	1709997/739032

6 MORE DETAILS ABOUT THE EXPERIMENTS

In this section we describe more detailed settings about the experiments to help in reproducibility.

6.1 Datasets and software versions

We describe more details about the datasets in Table 12. We download the datasets PPI, Reddit from the website⁸ and Amazon from the website⁹. Note that for Amazon, we consider GCN in an inductive setting, meaning that the model only learns from training data. In [3] they consider a transductive setting. Regarding software versions, we install CUDA 10.0 and cuDNN 7.0. TensorFlow 1.12.0 and PyTorch 1.0.0 are used. We download METIS 5.1.0 via the official website¹⁰ and use a Python wrapper¹¹ for METIS library.

6.2 Implementation details

Previous works [1, 2] propose to pre-compute the multiplication of AX in the first GCN layer. We also adopt this strategy in our implementation. By precomputing AX , we are essentially using the exact 1-hop neighborhood for each node and the expensive neighbors searching in the first layer can be saved.

Another implementation detail is about the technique mentioned in Section 3.2 When multiple clusters are selected, some between-cluster links are added back. Thus the new combined adjacency matrix should be re-normalized to maintain numerical ranges of the resulting embedding matrix. From experiments we find the renormalization is helpful.

As for the inductive setting, the testing nodes are not visible during the training process. Thus we construct an adjacency matrix containing only training nodes and another one containing all nodes. Graph partitioning are applied to the former one and the partitioned adjacency matrix is then re-normalized. Note that feature normalization is also conducted. To calculate the memory usage, we consider `tf.contrib.memory_stats.BytesInUse()` for TensorFlow and `torch.cuda.memory_allocated()` for PyTorch.

6.3 The running time of graph clustering algorithm and data preprocessing

The experiments of comparing different GCN training methods in Section 4 consider running time for training. The preprocessing time for each method is not presented in the tables and figures.

Table 13: The running time of graph clustering algorithm (METIS) and data preprocessing before the training of GCN.

Datasets	#Partitions	Clustering	Preprocessing
PPI	50	1.6s	20.3s
Reddit	1500	33s	286s
Amazon	200	0.3s	67.5s
Amazon2M	15000	148s	2160s

While some of these preprocessing steps such as data loading or parsing are shared across different methods, some steps are algorithm specific. For instance, our method needs to run graph clustering algorithm during the preprocessing stage.

In Table 13, we present more details about preprocessing time of Cluster-GCN on the four GCN datasets. For graph clustering, we adopt Metis, which is a fast and scalable graph clustering library. We observe that the graph clustering algorithm only takes a small portion of preprocessing time, showing a small extra cost while applying such algorithms and its scalability on large data sets. In addition, graph clustering only needs to be conducted once to form the node partitions, which can be re-used for later training processes.

⁸<http://snap.stanford.edu/graphsage/>

⁹https://github.com/HanJun-Dai/steady_state_embedding

¹⁰<http://glaros.dtc.umn.edu/gkhome/metis/metis/download>

¹¹<https://metis.readthedocs.io/en/latest/>