

SUGAR: Efficient Subgraph-level Training via Resource-aware Graph Partitioning

Zihui Xue, Yuedong Yang, and Radu Marculescu, *Fellow, IEEE*

Abstract—Graph Neural Networks (GNNs) have demonstrated a great potential in a variety of graph-based applications, such as recommender systems, drug discovery, and object recognition. Nevertheless, resource-efficient GNN learning is a rarely explored topic despite its many benefits for edge computing and Internet of Things (IoT) applications. To improve this state of affairs, this work proposes efficient subgraph-level training via resource-aware graph partitioning (SUGAR). SUGAR first partitions the initial graph into a set of disjoint subgraphs and then performs local training at the subgraph-level. We provide a theoretical analysis and conduct extensive experiments on five graph benchmarks to verify its efficacy in practice. Our results across five different hardware platforms demonstrate great runtime speedup and memory reduction of SUGAR on large-scale graphs. We believe SUGAR opens a new research direction towards developing GNN methods that are resource-efficient, hence suitable for IoT deployment. Our code is publicly available at: <https://github.com/zihuixue/SUGAR>.

Index Terms—Graph Neural Networks, Resource-efficient Learning, Edge Computing

I. INTRODUCTION

GRAPHS are non-Euclidean data structures that can model complex relationships among a set of interacting objects, for instance, social networks, knowledge graphs, or biological networks. Given the huge success of deep neural networks for Euclidean data (*e.g.*, images, text and audio), there is an increasing interest in developing deep learning approaches for graphs too. Graph Neural Networks (GNNs) generalize the convolution operation to the non-Euclidean domain [1]; they demonstrate a great potential for various graph-based applications, such as node classification [2], link prediction [3] and recommender systems [4].

The rapid development of smart devices and IoT applications has spawned a great interest in

many edge AI applications. Training models locally becomes a growing trend as this can help avoid data transmission to the cloud, reduce communication latency, and better preserve privacy [5]. For instance, in a graph-based recommender system, user data can be quite sensitive and hence it's better to store it locally [6]. This brings about the need for *resource-efficient graph learning*.

While there is much discussion about locally training Convolutional Neural Networks (CNNs) [7], efficient on-device training for GNNs is rarely explored. Different from CNNs, where popular models such as ResNet [8] are deep and have a large parameter space, mainstream GNN models are shallow and more lightweight. However, the major bottleneck of GNN training comes from the nodes dependencies in the input graph. Consequently, graph convolution suffers from a high computational cost, as the representation of a node in the current layer needs to be computed recursively by the representations of all neighbors in its previous layer. Moreover, storing the intermediate features for all nodes requires much memory space, especially when the graph size grows. For instance, for the *ogbn-products* graph in our experiments (Table I), full-batch training requires a GPU with 33GB of memory [9]. Thus scaling GNN training to large-scale graphs remains a big challenge. The problem is more severe for a resource-constrained scenario like IoT, where GNN training is heavily constrained by the computation, memory, and communication costs.

Various approaches have been proposed to alleviate the computation and memory burden of GNNs. For instance, sampling-based approaches aim at reducing the neighborhood size via layer sampling [10]–[12], clustering based sampling [13] and graph sampling [14] techniques; these prior works approach this problem purely from an algorithmic angle. A few recent works [15], [16] investigate the topic of distributed multi-GPU training of GNNs and achieve good parallel efficiency and memory scalability while

Zihui Xue, Yuedong Yang and Radu Marculescu are with the Department of Electrical and Computer Engineering at The University of Texas at Austin.
E-mail: {sherryxue, albertyoung, radum}@utexas.edu

using large GPU clusters.

A common limitation of all these approaches is that they *do not* take the real hardware constraints into consideration. For mobile devices with limited memory budgets, the input graph can be too large to fit entirely in the main memory. In addition, the communication overhead among real IoT devices is significantly larger than when using GPU clusters, rendering distributed training approaches not readily applicable to such scenarios. This calls for a new approach for *resource-efficient GNN learning*, which is precisely the focus of our paper.

In this work, we propose a novel approach that trains GNNs efficiently with multiple devices in a resource-limited scenario. To be specific, we assume that only several resource-constrained devices are available for GNN training, and no inter-device communication is allowed. To this end, we (1) design a graph partitioning method that accounts for resource constraints and graph topology; (2) train a set of local GNNs at the *subgraph-level* for computation, memory and communication savings. Our contributions are as follows:

- We formulate the problem of training GNNs with multiple resource-constrained devices. Although our formulation targets various mobile and edge devices (*e.g.*, mobile phones, Raspberry Pi), it is also applicable to powerful machines equipped with GPUs.
- We propose SUGAR, a GNN training framework that aims at improving training scalability. We provide complexity analysis, error bound and convergence analysis of the proposed estimator.
- We show that SUGAR achieves the best runtime and memory usage (with similar accuracy) when compared against state-of-the-art GNN approaches on five large-scale datasets and across multiple hardware platforms, ranging from edge devices (*i.e.*, Raspberry Pi, Jetson Nano) to a desktop equipped with powerful GPUs.
- We illustrate the flexibility of SUGAR by integrating it with both full-batch and mini-batch algorithms such as GraphSAGE [10] and GraphSAINT [14]. Experimental results demonstrate that SUGAR can achieve up to

33× runtime speedup on *ogbn-arxiv* and 3.8× memory reduction on *Reddit*. On the *ogbn-products* graph with over 2 million nodes and 61 million edges, SUGAR achieves 1.62× speedup over GraphSAGE and 1.83× memory reduction over GraphSAINT with a better test accuracy ($\sim 0.7\%$).

The remainder of the paper is organized as follows. In Section 2, we discuss prior work. In Section 3, we formulate the problem and describe our proposed training framework SUGAR. Experimental results are presented in Section 4. Finally, Section 5 concludes the paper.

II. RELATED WORK

The relevant prior work comes from four directions as discussed next.

A. Graph Neural Networks

Modern GNNs adopt a neighborhood aggregation scheme to learn representations for individual nodes or the entire graph. Graph Convolution Network (GCN) [2] is a pioneering work that generalizes the use of regular convolutions to graphs. GraphSAGE [10] provides an inductive graph representation learning framework. To improve the representation ability of GNNs, Graph Attention Networks (GAT) [17] introduce self-attention to the graph convolution operation. Apart from pursuing higher accuracy, a few GNN architecture improvements [18], [19] have been made towards higher training efficiency.

B. GNN Training Algorithms

Full-batch training was first proposed for GCNs [2]; the gradient is calculated based on the global graph and updated once per epoch. Despite being fast, full-batch gradient descent is generally infeasible for large-scale graphs due to excessively large memory requirements and slow convergence.

Mini-batch training was first proposed in GraphSAGE [10]; the gradient update is based on a proportion of nodes in the graph and updated a few times during each training epoch. Mini-batch training leads to memory efficiency at the cost of increased computation. Since the neighborhood aggregation scheme involves recursive calculation of a node's neighbors layer by layer, time complexity becomes exponential with respect to the number of GNN

layers; this is known as the *neighborhood expansion problem*. Following the idea of neighbor sampling, FastGCN [11] further proposes the importance node sampling to reduce variance. The work of [12] proposes a control variate based algorithm that allows a smaller neighbor sample size. A few recent works propose alternative ways to construct mini-batches instead of layer-wise sampling. For instance, ClusterGCN [13] first partitions the training graph into clusters and then randomly groups clusters together as a batch. GraphSAINT [14] builds mini-batches by sampling the training graph and ensures a fixed number of nodes in all layers.

Distributed training aims at leveraging multiple devices to speed up the training process. Many distributed GNN training approaches [15], [16], [20] have been developed to enable the exploration towards larger models and datasets. DistDGL [16] distributes the input graph across machines via METIS-based graph partitioning and performs synchronous training. GIST [15] proposes to partition the parameters of a GCN model into smaller sub-GCNs and train several sub-GCNs in parallel. The recent work BNS-GCN [20] proposes random boundary node sampling to enable efficient and scalable distributed training. While previous works generally assume powerful GPU clusters and allow inter-device communication, in this work, we target on IoT scenarios with only resource-constrained edge devices available and no communication.

C. Graph Partitioning

Graph partitioning is a widely studied topic in the field of graph processing. A variety of methods have been proposed to divide a large graphs into a predefined number of subgraphs. One popular approach is METIS [21], which recursively partitions the graph via k-way partitioning and recursive bisection. METIS serves as a powerful tool for distributed graph processing and parallel computing due to its ability to achieve a good load balance and minimal communication cost. The effectiveness of METIS in partitioning large graphs makes it a widely used method in modern GNN training.

The growing need for low-latency, continuous graph analysis has led to the development of online partitioning methods [22], which ingest the graph data as a stream and making partitioning decisions on the fly based on partial knowledge of the graph.

Among them, HDRF [23] and DBH [24] exploit skewed degree distributions in power-law graphs and explicitly account for vertex degree in the placement decision. In this work, we also utilize vertex degree information to facilitate the partitioning process, yet the input graph is not restricted to power-law graphs.

D. Graph Sparsification

Recent works have also investigated graph sparsification (*i.e.*, pruning edges of the training graph) for GNN learning. In many real-world applications, graphs exhibit complex topology patterns. Some edges may be erroneous or task-irrelevant, and thus aggregating this information weakens the generalizability of GNNs [25]. As shown by [26] and [27], edges of the input graph may be pruned without loss of accuracy.

Two recent works introduce computation efficiency into the problem. More precisely, SGCN [28] proposes a neural network that prunes edges of the input graph; they show that using sparsified graphs as the new input for GNNs brings computational benefits. UGS [29] presents a graph lottery ticket type of approach; they sparsify the input graph, as well as model weights during training to save inference computation.

III. OUR PROPOSED METHOD

A. Problem Formulation

Given a graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, where \mathcal{V} is the node set and \mathcal{E} represents the set of edges. Let $N = |\mathcal{V}|$ denote the number of nodes and $A \in R^{N \times N}$ be the adjacency matrix of \mathcal{G} . Every node i is characterized by a F -dimensional feature vector $x_i \in R^F$. We use $X \in R^{N \times F}$ to represent the feature matrix of all nodes in \mathcal{G} .

Consider a node-level prediction problem with the following objective:

$$\min_W \mathcal{L} = \frac{1}{N} \sum_{i=1}^N f(y_i, z_i) \quad (1)$$

$$z_i = g(x_i; W)$$

where f is the objective function (*e.g.*, cross entropy for node classification), y_i and z_i denotes the true label and prediction of node i , respectively. $g(\cdot)$ denotes a graph neural network parameterized by W that generates node-level predictions.

Suppose there are K devices available for training, and let \mathcal{B}_{MEM}^k denote the memory budget of device k . Motivated by the notorious inefficiency that centralized graph learning suffers from, we aim at *distributing the training process* to improve the training scalability with no inter-device communication overhead. The key is to assign N nodes of graph \mathcal{G} to K devices, and then do local training on *each* device. We formulate it as two subproblems below.

First, we define a graph partitioning strategy $\mathcal{P} : \mathcal{V} \rightarrow (\mathcal{V}_1, \mathcal{V}_2, \dots, \mathcal{V}_K)$ that divides the node set \mathcal{V} into K subsets such that:

$$\cup_k \mathcal{V}_k = \mathcal{V}, \quad H(\mathcal{SG}_k) < \mathcal{B}_{MEM}^k, \quad \forall k \in [K] \quad (2)$$

where $[K] = \{1, \dots, K\}$, \mathcal{SG}_k denotes the subgraph induced by node set \mathcal{V}_k , H is a static function that maps a given subgraph \mathcal{SG}_i to the device memory requirements for training. For maximum generality, here we do not require $\mathcal{V}_i \cap \mathcal{V}_j = \emptyset$. In other words, a node i can be assigned to more than one hardware device, and let \mathcal{P}_i denote the set of hardware devices where node i is assigned to.

Next, we adopt subgraph-level training, *i.e.*, for device k , we maintain a local GNN model, denoted by $W^{(k)}$ that takes the subgraph \mathcal{SG}_k as its input graph. Let $W = \frac{1}{K} \sum_{k=1}^K W^{(k)}$, thus the objective can be reformulated as:

$$\begin{aligned} \min_W \mathcal{L} &= \frac{1}{N} \sum_{i=1}^N f(y_i, z_i) \\ z_i &= \frac{1}{|\mathcal{P}_i|} \sum_{k \in \mathcal{P}_i} g(x_i; W^{(k)}) \end{aligned} \quad (3)$$

Based on the formulation above, we propose SUGAR, a distributed training framework that: (1) partitions the input graph subject to resource constraints; (2) adopts local subgraph-level training. Figure 1 provides a simple illustration of SUGAR for a two-device system. We describe our design choices in detail in the following sections.

B. Theoretical Basis

Recall that we define a graph partitioning strategy \mathcal{P} that divides N nodes into K node sets $(\mathcal{V}_1, \mathcal{V}_2, \dots, \mathcal{V}_K)$. Taking K subgraphs induced by the node sets into consideration, a graph partitioning strategy \mathcal{P} can be viewed as a way to produce a sparser adjacency matrix A_{SG} , from the original

matrix A . A_{SG} is a block-diagonal matrix of A , *i.e.*,

$$A_{SG} = \begin{bmatrix} A_{\mathcal{V}_1} & \cdots & 0 & \cdots & 0 \\ \vdots & \ddots & & & \vdots \\ 0 & & A_{\mathcal{V}_k} & & 0 \\ \vdots & & & \ddots & \vdots \\ 0 & \cdots & 0 & \cdots & A_{\mathcal{V}_K} \end{bmatrix} \quad (4)$$

where $A_{\mathcal{V}_k}$ denotes the adjacency matrix of subgraph k .

We show below that adopting A_{SG} for training offers the benefits of high computational efficiency and low memory requirements. Moreover, we provide the error bound and convergence analysis of this approximation for a graph convolutional network (GCN) [2].

Complexity Analysis. The propagation rule for the l -th layer GCN is:

$$Z^{(l+1)} = A^{norm} H^{(l)} W^{(l)}, \quad H^{(l+1)} = \sigma(Z^{(l+1)}) \quad (5)$$

where σ represents an activation function, A^{norm} denotes the normalized version of A , *i.e.*, $A^{norm} = \hat{D}^{-1/2} \hat{A} \hat{D}^{1/2}$, $\hat{A} = A + I_N$, $\hat{D}_{ii} = \sum_j \hat{A}_{ij}$ and I_N is an N -dimensional identity matrix. $H^{(l)}$ and $H^{(l+1)}$ denotes the input and output feature matrices in layer l , respectively. $Z^{(l)}$ is the node feature matrix before the activation function in layer l and $Z^{(L)}$ denotes final node predictions (*i.e.*, output of the GCN). $W^{(l)} \in R^{F_l \times F_{l+1}}$ represents the weight matrix of layer l , where F_l and F_{l+1} is the input and output feature dimension, respectively. Therefore, for the l -th layer GCN, the *training time complexity* is $\mathcal{O}(|\mathcal{E}| F_l + N F_l F_{l+1})$ and *memory complexity* is $\mathcal{O}(N F_{l+1} + F_l F_{l+1})$. We make two observations here: (a) Real-world graphs are usually sparse and $\frac{|\mathcal{E}|}{N}$ is generally smaller than feature number F_{l+1} . Thus, the second term dominates the time complexity; (b) For large-scale graphs, the number of nodes N is much greater than the number of features, so $\mathcal{O}(N F_{l+1})$ dominates the memory complexity. To conclude, the value of N affects the complexity greatly; indeed, calculating these terms are likely to result in poor data locality and dominate the runtime latency due to the high randomness of neighbor indices [30]. Partitioning the input graph into K subgraphs reduces the number of nodes N to $N_k = |\mathcal{V}_k|$ for every local model. Since N_k is about $1/K$ of N , the proposed approach is expected to achieve up to K times speedup, and as little as $1/K$ of the original memory requirements.

Error Bound Analysis. Let our proposed estimator be SG. The l -th layer propagation rule of a GCN with the SG estimator is:

$$Z_{SG}^{(l+1)} = A_{SG}^{norm} H_{SG}^{(l)} W^{(l)}, H_{SG}^{(l+1)} = \sigma(Z_{SG}^{(l+1)}) \quad (6)$$

where $Z_{SG}^{(l+1)}$ and $H_{SG}^{(l+1)}$ denote the node representations produced by the SG estimator in layer $l+1$ before and after activation, respectively.

Assume that we run graph partitioning for M times to obtain a sample average of A_{SG}^{norm} before training. Let $\epsilon = \|A_{SG}^{norm} - A^{norm}\|_\infty$ denote the error in approximating A^{norm} with A_{SG}^{norm} . For simplicity, we will omit the superscript *norm* from now on.

The following lemma states that the error of node predictions given by the SG estimator is bounded.

Lemma 1. *For a multi-layer GCN with fixed weights, assume that: (1) $\sigma(\cdot)$ is ρ -Lipschitz and $\sigma(0) = 0$, (2) input matrices A , X and model weights $\{W^{(l)}\}_{l=1}^L$ are all bounded, then there exists C such that $\|Z_{SG}^{(l)} - Z^{(l)}\|_\infty \leq C\epsilon, \forall l \in [L]$ and $\|H_{SG}^{(l)} - H^{(l)}\|_\infty \leq C\epsilon, \forall l \in [L-1]$.*

The proof of Lemma 1 is provided in Appendix. Lemma 1 motivates us to design a graph partitioning method that generates small ϵ so that the output of the SG estimator is close to the exact value. This will be discussed in detail in the next subsection.

Convergence Analysis. Let W_t denote the model parameters at training epoch t and W_* denote the optimal model weights. $\nabla \mathcal{L}(W) = \frac{1}{N} \sum_{i=1}^N \frac{\partial f(y_i, z_i^{(L)})}{\partial W}$ and $\nabla \mathcal{L}_{SG}(W) = \frac{1}{N} \sum_{i=1}^N \frac{\partial f(y_i, z_{SG,i}^{(L)})}{\partial W}$ represent the gradients of the exact GCN and SG estimator with respect to model weights W , respectively.

Theorem 1 states that with high probability gradient descent training with the approximated gradients of the SG estimator (i.e., $\nabla \mathcal{L}_{SG}(W)$) converges to a local minimum.

Theorem 1. *Assume that: (1) the loss function $\mathcal{L}(W)$ is ρ -smooth, (2) the gradients of the loss $\nabla \mathcal{L}(W)$ and $\nabla \mathcal{L}_{SG}(W)$ are bounded for any choice of W , (3) the gradient of the objective function $\frac{\partial f(y,z)}{\partial z}$ is ρ -Lipschitz and bounded, (4) the activation function $\sigma(\cdot)$ is ρ -Lipschitz, $\sigma(0) = 0$ and its gradient is bounded,*

then there exists $C > 0$, s.t., $\forall M, T$, for a sufficiently small δ , if we run graph partitioning

for M times and run gradient descent for $R \leq T$ epochs (where R is chosen uniformly from $[T]$, the model update rule is $W_{t+1} = W_t - \gamma \nabla \mathcal{L}_{SG}(W_t)$, and step size $\gamma = \frac{1}{\rho\sqrt{T}}$), we have:

$$P(\mathbb{E}_R \|\nabla \mathcal{L}(W_R)\|_F^2 \leq \delta) \geq 1 - 2 \exp\left\{-2M\left(\frac{\delta}{2C} - \frac{2\rho[\mathcal{L}(W_1) - \mathcal{L}(W_*)] + C - \delta}{2C(\sqrt{T} - 1)}\right)^2\right\}$$

With M and T increasing, the right-hand-side of the inequality becomes larger. This implies that there is a higher probability for the loss to converge to a local minimum. The full proof is provided in the Appendix.

C. Graph Partitioning

From Lemma 1, we conclude that a graph partitioning method that yields a smaller $|A_{SG} - A|$ leads to a smaller error in node predictions. Therefore, we aim at minimizing the difference between A_{SG} and A . In other words, the objective of graph partitioning should be to *minimize the number of edges of the incident nodes* that belong to different subsets. As such, this is identical to the goal of various existing graph partitioning methods, making such approaches good candidates to use with our framework. We choose METIS [21] due to its efficiency in handling large-scale graphs. However, the traditional graph partitioning algorithms are *not* intended for modern GNNs and the learning component of the problem is missing. Consequently, we present a modified version of METIS that is suited to our problem and relies on two new ideas discussed next.

a) Weighted Graph Construction. We build a weighted graph \mathcal{G}^w from the input graph \mathcal{G} . The *weight* of an edge e_{uv} is defined based on the degree of its two incident nodes:

$$\begin{aligned} \text{weight}(e_{uv}) &= d_{max} + 1 - \deg(u) - \deg(v) \\ d_{max} &= \max\{\deg(u) + \deg(v), \forall e_{uv} \in \mathcal{E}\} \end{aligned} \quad (7)$$

Let A^w denote the adjacency matrix of the weighted graph \mathcal{G}^w , where element a_{ij}^w is the edge weight $\text{weight}(e_{ij})$; a_{ij}^w is 0 if there is no edge connecting nodes i and j .

The key intuition behind our first idea lies in the neighborhood aggregation scheme of GNNs. Consider two nodes u and v , where u is a hub node connected to many other nodes, while v has only one neighbor. As GNNs propagate by aggregating

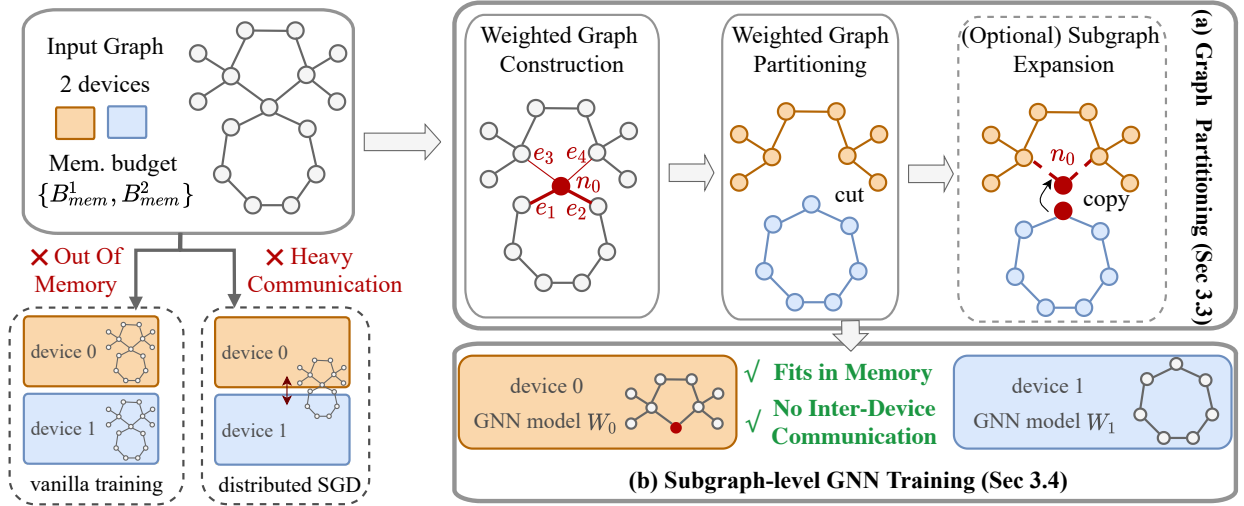


Fig. 1. While vanilla training is likely to run out of memory when the graph size is large and distributed stochastic gradient descent (SGD) requires heavy intermediate communication among devices, SUGAR provides a solution that is memory efficient and requires no inter-device communication. The proposed SUGAR consists of two stages: (a) graph partitioning and (b) subgraph-level GNN training. Graph partitioning involves three steps: (1) transform the input graph \mathcal{G} to a weighted graph \mathcal{G}^w ; (2) apply METIS to the weighted graph \mathcal{G}^w , where edges with large weights are more likely to be preserved; (3) (optional) expand the node set of the obtained subgraph according to memory budgets.

the neighborhood information of nodes, removing the only edge of node v may possibly lead to wrong predictions. On the other hand, pruning an edge of u is more acceptable since there are many neighbors contributing to its prediction. Consider the graph in Figure 1 as an example. Cutting the edges $e_1 \cup e_2$ and $e_3 \cup e_4$ are both feasible solutions for METIS. However, considering the fact that nodes connected to e_1 and e_2 have less topology information, our proposed method will preserve them and cut edges $e_3 \cup e_4$ instead; this can lead to a better learning performance.

As can be concluded from this small example, edges connected to small-degree nodes are critical to our problem and should be preserved. Conversely, edges connected to high-degree nodes may be intentionally ignored. This explains our weights definition strategy. Consequently, we incorporate the above observation into our partitioning objective and apply METIS to the pre-processed graph \mathcal{G}^w .

b) Subgraph Expansion. After obtaining the partitions with our modified METIS, we propose the second idea, *i.e.*, expand the subgraph based on available hardware resources. Although METIS only provides partitioning results where the node sets do not overlap, our general formulation in Section III-A allows nodes to belong to multiple partitions. This brings great flexibility to our approach to adjust the node number for each device according to its

memory budget.

Suppose the available memory of device k is larger than the actual requirement of training a GNN on subgraph k (*i.e.*, $H(\mathcal{SG}_k) < B_{MEM}^k$), then we may choose to expand the node set \mathcal{V}_k by adding the one-hop neighbors of nodes that do not belong to \mathcal{V}_k . As illustrated in Figure 1 (a), we can expand the node set of the subgraph on device 0 (marked in light brown) to include node n_0 as well. While expanding the subgraph is likely to yield higher accuracy, training time and memory requirement will also increase. Therefore, this is an optional step, only if the hardware resources allow it.

D. Subgraph-level Local Training

From the original formulation in Equation 3, if $|\mathcal{P}_i| > 1$, *i.e.*, a node i is assigned to multiple devices, calculating its loss and backpropagation can involve heavy communication among devices. To address this problem, we provide the following result to decouple the training of K local GNN models from each other.

Proposition 1. *If $f(y, z)$ is convex with respect to z , then the upper bound of \mathcal{L} in Equation 3 is given by:*

$$\mathcal{L} \leq \frac{1}{K} \sum_{k=1}^K \sum_{i \in \mathcal{V}_k} \frac{1}{|\mathcal{P}_i|} f(y_i, z_i) \quad (8)$$

$$z_i = g(x_i, W^{(k)})$$

Proof. By convexity of f , using Jensen's inequality [31] gives us:

$$f(y_i, \frac{1}{|\mathcal{P}_i|} \sum_{k \in \mathcal{P}_i} g(x_i; W^{(k)}) \leq \frac{1}{|\mathcal{P}_i|} \sum_{k \in \mathcal{P}_i} f(y_i, g(x_i; W^{(k)})) \quad (9)$$

By changing the operation order and regrouping the indices, we further derive:

$$\frac{1}{N} \sum_{i=1}^N \frac{1}{|\mathcal{P}_i|} \sum_{k \in \mathcal{P}_i} f(y_i, z_i) = \frac{1}{K} \sum_{k=1}^K \sum_{i \in \mathcal{V}_k} \frac{1}{|\mathcal{P}_i|} f(y_i, z_i) \quad (10)$$

Therefore,

$$\mathcal{L} \leq \frac{1}{K} \sum_{k=1}^K \sum_{i \in \mathcal{V}_k} \frac{1}{|\mathcal{P}_i|} f(y_i, z_i) \quad (11)$$

$$z_i = g(x_i, W^{(k)})$$

Proposition 1 is proved.

Proposition 1 allows us to shift the perspective from 'node-level' to 'device-level'. We adopt the upper bound of \mathcal{L} in Equation 8 as the new training objective. Now, the local model updates involving node i do not depend on other models (*i.e.*, $\{W^{(k)}\}_{k \in \mathcal{P}_i}$) any more. Optimizing the new objective naturally reduces the upper bound of the original one and avoids significant communication costs, thus leading to high training efficiency.

Furthermore, motivated by deployment challenges in real IoT applications, where communication among devices is generally not guaranteed, we propose to reduce inter-device communication down to zero in our framework. In particular, we maintain K distinct (local) models instead of a single (global) model by keeping the local model updates within each device. The objective of our proposed subgraph-level local GNN training can be summarized as follows:

$$\min_{W^{(k)}} \mathcal{L}_k = \sum_{i \in \mathcal{V}_k} \frac{1}{|\mathcal{P}_i|} f(y_i, z_i), \quad \forall k \in [K] \quad (12)$$

$$z_i = g(x_i, W^{(k)})$$

In training round t , every device performs local updates as:

$$W_{t+1}^{(k)} \leftarrow W_t^{(k)} - \gamma \nabla_{W^{(k)}} \mathcal{L}_k, \quad \forall k \in [K] \quad (13)$$

where \mathcal{L}_k denotes the training objective of device k and γ is the learning rate (*i.e.*, step size). By

decoupling training dependency among devices, we propose a feasible solution to train GNNs in resource-limited scenarios, where typical distributed GNN approaches are not applicable.

E. Putting it all together

Algorithm 1 SUGAR

Input: graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$; node feature matrix X ; available device number K ; device memory budget $\{B_{MEM}^k\}_{k=1}^K$; total training epochs T .

- 1: Construct \mathcal{G}^w from \mathcal{G} according to Equation 7
 - 2: Partition \mathcal{G}^w into K subgraphs $\{\mathcal{SG}_i\}_1^K$
 - 3: (Optional) Expand \mathcal{SG}_i if $H(\mathcal{SG}_i) < B_{MEM}^i$
 - 4: **for** each device $k = \{1, 2, \dots, K\}$ in parallel **do**
 - 5: Initialize GNN model weight $W_1^{(k)}$
 - 6: **for** epoch $t = 1, 2, \dots, T$ **do**
 - 7: $W_{t+1}^{(k)} \leftarrow W_t^{(k)} - \gamma \nabla_{W^{(k)}} \mathcal{L}_k$
 - 8: **end for**
 - 9: **end for**
-

To sum up, the SUGAR algorithm consists of two stages: (a) graph partitioning (lines 1-3) and (b) subgraph-level GNN training (lines 4-9). Specifically, the graph partitioning involves three steps: (1) construct a weighted graph \mathcal{G}^w from \mathcal{G} to account for the influence of node degrees in learning (line 1). (2) Apply METIS to the weighted graph \mathcal{G}^w to obtain partitioning results (line 2). (3) According to the memory budget, expand the subgraph to cover the one-hop neighbors for better performance (line 3). Then, we train K local models in parallel without requiring training-time communication among devices (lines 4-9). During inference, we adopt the local GNN model $W^{(k)}$ to generate predictions for test nodes that belong to \mathcal{SG}_k . Similar to training, the inference process is distributed to K devices and requires no inter-device communication. In conclusion, SUGAR provides a highly efficient and cost-effective solution in scaling GNN training, with high training speed, low memory requirements and no communication overhead.

IV. EXPERIMENTS

A. Experimental Setup

We evaluate SUGAR on five node classification datasets [9], [32], selected from very diverse applications: (1) categorizing types of images based

TABLE I

DATASET STATISTICS. K AND M DENOTE 1,000 AND 1,000,000, RESPECTIVELY. ‘AVGDEG.’ REPRESENTS THE AVERAGE NODE DEGREE. ‘ACC’ DENOTES ACCURACY. ‘ROC-AUC’ DENOTES THE AREA UNDER THE RECEIVER OPERATING CHARACTERISTIC CURVE.

Dataset	<i>Flickr</i>	<i>Reddit</i>	<i>ogbn-arxiv</i>	<i>ogbn-proteins</i>	<i>ogbn-products</i>
#Nodes	89.3K	233K	169K	133K	2,449K
#Edges	0.90M	11.6M	1.17M	39.6M	61.9M
AvgDeg.	10	50	13.77	597	50.5
#Tasks	1	1	1	112	1
#Classes	7	41	40	2	47
Metric	ACC	ACC	ACC	ROC-AUC	ACC

on the descriptions and common properties of online images (*Flickr*); (2) predicting communities of online posts based on user comments (*Reddit*); (3) predicting the subject areas of arxiv papers based on its title and abstract (*ogbn-arxiv*); (4) predicting the presence of protein functions based on biological associations between proteins (*ogbn-proteins*); (5) predicting the category of a product in an Amazon product co-purchasing network (*ogbn-products*). Note that the task of *ogbn-proteins* is multi-label classification, while other tasks are multi-class classification. Dataset statistics are summarized in Table I.

We include the following GNN architectures and SOTA GNN training algorithms for comparison:

- GCN [2]: Full-batch Graph Convolutional Networks.
- GraphSAGE [10]: An inductive representation learning framework that efficiently generates node embeddings for previously unseen data. Mini-batch GraphSAGE are denoted by GraphSAGE-mb.
- GAT [17]: Graph Attention Networks, a GNN architecture that leverages masked self-attention layers.
- SIGN [18]: Scalable Inception Graph Neural Networks, an architecture using graph convolution filters of different size for efficient computation.
- ClusterGCN [13]: A mini-batch training technique that partitions the graphs into a fixed

TABLE II

RUNTIME, MEMORY & ACCURACY RESULTS ON *ogbn-arxiv*. ‘AVG. TIME’ IS THE TRAINING TIME PER EPOCH AVERAGED OVER 100 EPOCHS AND ‘MAX MEM’ DENOTES PEAK ALLOCATED MEMORY ON GPU.

	Avg. Time [ms]	SUGAR Speedup	Max Mem [GB]	Test Acc. [%]
GCN	26.9	1.68×	1.60	72.37 ± 0.10
GAT	207.8	12.99×	5.41	72.95 ± 0.14
GraphSAGE	534.7	33.42×	0.95	71.98 ± 0.17
SIGN	291.6	18.23×	0.94	71.79 ± 0.08
SUGAR	16.0	---	0.92	72.22 ± 0.14

number of subgraphs and draws mini-batches from them.

- GraphSAINT [14]: A mini-batch training technique that constructs mini-batches by graph sampling. The random node, random edge, and random walk based samplers are denoted by GraphSAINT-N, GraphSAINT-E, GraphSAINT-RW, respectively.

SUGAR is implemented with PyTorch [33] and DGL [34]. For all the baseline methods, we use the parameters reported in their github pages or the original paper. The evaluation of SUGAR is conducted on a two-device system (*i.e.*, $K = 2$) unless otherwise stated. In Section IV-C1, we provide a scalability analysis with varying K . The baseline GNN approaches are evaluated on a single device as we assume no inter-device communication in our setting. For the dense *ogbn-proteins* graph, after graph partitioning, we expand the node set to include an additional 20% of the nodes in the subgraph. For all other graphs, we do not adopt the subgraph expansion step. We report accuracy results averaged over 5 runs for *ogbn-proteins* and 10 runs for the other datasets.

For completeness, we run our experiments across multiple hardware platforms. We select five different devices with various computing and memory capabilities, namely, (1) Raspberry Pi 3B, (2) NVIDIA Jetson Nano, (3) Android phone with Snapdragon 845 processor, (4) laptop with Intel i5-8279U CPU, and (5) desktop with AMD Threadripper 3970X CPU and two NVIDIA RTX 3090 GPUs.

TABLE III
RUNTIME, MEMORY & ACCURACY RESULTS ON *Reddit*.

	Avg. Time [ms]	SUGAR Speedup	Max Mem [GB]	Test Acc. [%]
GraphSAGE	110.6	1.87×	5.70	96.39 ± 0.03
GraphSAGE-mb	316.5	5.36×	2.33	95.08 ± 0.05
ClusterGCN	414.4	7.01×	1.83	96.34 ± 0.01
GraphSAINT-N	341.8	5.78×	1.29	96.17 ± 0.06
GraphSAINT-E	299.8	5.07×	1.22	96.15 ± 0.06
GraphSAINT-RW	467.5	7.91×	1.23	96.23 ± 0.06
SIGN	352.8	5.97×	2.17	96.12 ± 0.05
SUGAR	59.1	1.51	1.51	96.01 ± 0.03

TABLE IV
RUNTIME, MEMORY & ACCURACY RESULTS ON *Flickr*.

	Avg. Time [ms]	Max Mem [GB]	Test Acc. [%]
GraphSAINT-N	97.0	0.41	50.64 ± 0.28
SUGAR	49.9	0.31	50.11 ± 0.12
Improvement	1.94×	1.32×	
GraphSAINT-E	71.1	0.53	50.91 ± 0.12
SUGAR	32.6	0.41	49.96 ± 0.12
Improvement	2.18×	1.29×	
GraphSAINT-RW	108.9	0.65	51.03 ± 0.20
SUGAR	37.3	0.49	50.15 ± 0.24
Improvement	2.92×	1.33×	

B. Results

1) *Evaluations on GPUs*: First, we provide evaluation of SUGAR on a two-GPU system. Table II and Table III report the average training time per epoch, maximum GPU memory usage and accuracy on *ogbn-arxiv* and *Reddit*. We base SUGAR on full-batch GCN and GraphSAGE for these two datasets, respectively. As shown in these tables, when compared with full-batch methods (*i.e.*, GCN and GAT for *ogbn-arxiv*; GraphSAGE for *Reddit*), SUGAR is much more memory efficient, as it reduces the peak memory by 1.7× for *ogbn-arxiv* and 3.8× for *Reddit* data. When compared against mini-batch methods (*i.e.*, mini-batch GraphSAGE, ClusterGCN, GraphSAINT and SIGN), the runtime of SUGAR is significantly smaller. This demonstrates the great benefits of our proposed subgraph-level training. Indeed, by restricting the neighborhood search size, SUGAR effectively alleviates the neighborhood expansion problem. In addition, it achieves very competitive test accuracies.

We combine SUGAR with popular mini-batch

TABLE V
RUNTIME, MEMORY & ACCURACY RESULTS ON *ogbn-products*.

	Avg. Time [ms]	Max Mem [GB]	Test Acc. [%]
GraphSAGE-mb	2.42	7.29	79.25 ± 0.22
SUGAR	1.49	4.43	79.97 ± 0.23
Improvement	1.62×	1.65×	
ClusterGCN	2.90	6.59	78.51 ± 0.33
SUGAR	1.97	3.36	79.34 ± 0.41
Improvement	1.47×	1.96×	
GraphSAINT-E	0.30	7.16	79.54 ± 0.27
SUGAR	0.28	3.92	80.20 ± 0.23
Improvement	1.07×	1.83×	

TABLE VI
RUNTIME, MEMORY & ACCURACY RESULTS ON *ogbn-proteins*.

	Avg. Time [sec]	Max Mem [GB]	Valid Acc. [%]	Test Acc. [%]
GAT	6.20	10.77	92.08 ± 0.08	87.20 ± 0.17
SUGAR	4.09	6.22	92.51 ± 0.08	86.41 ± 0.18
Improvement	1.52×	1.73×		

training methods and evaluate them on *Flickr* and *ogbn-products* dataset. Table IV presents results of SUGAR incorporated with GraphSAINT for three sampler modes (*i.e.*, node, edge, and random walk based samplers) on *Flickr* data. Note that the accuracy we obtain (about 50%) is consistent with results in [14]. SUGAR achieves more than 2× runtime speedup and requires less memory than GraphSAINT. Test accuracy loss is within 1% in all cases.

For the largest *ogbn-products* dataset, we implement SUGAR together with three competitive GNN baselines, namely GraphSAGE, ClusterGCN and GraphSAINT. The results are summarized in Table V. SUGAR provides a better solution that leads to runtime speedup, memory reduction and even a slightly increased test accuracy for all three methods. We hypothesize that the graph partitioning eliminates some task-irrelevant edges in the original graph, and thus leads to better generalization of GNNs.

Table VI provides results on the dense *ogbn-proteins* graph. When it comes to training GNNs on dense graphs, memory poses a significant challenge due to the neighborhood expansion problem. The results show that GAT suffers from considerable

TABLE VII

AVERAGE TRAINING TIME PER EPOCH [SEC] OF SUGAR COMPARED WITH GRAPHSAINT AND GCN ON *Flickr* AND *ogbn-arxiv* DATA. WE RECORD THE TRAINING TIME ON FIVE PLATFORMS WITH CPU MODELS LISTED. OOM DENOTES OUT OF MEMORY. WE NOTE THAT TRAINING A GCN ON RASPBERRY PI 3B IS INFEASIBLE SINCE IT EXCEEDS MEMORY, WHILE SUGAR STILL WORKS.

Dataset		RPi 3B Cortex-A53	Jetson Cortex-A57	Phone SDM-845	Laptop i5-8279U	Desktop-CPU Zen2 3970X	Desktop-GPU RTX3090
<i>Flickr</i>	GraphSAINT-N	104.1	16.86	7.67	2.86	1.48	0.097
	SUGAR	48.2	7.61	3.54	1.21	0.67	0.050
	Speedup	2.16×	2.22×	2.17×	2.36×	2.24×	1.94×
<i>ogbn-arxiv</i>	GCN	OOM	28.10	21.96	13.80	5.16	0.027
	SUGAR	501.59	18.39	13.33	6.51	2.71	0.016
	Speedup	-	1.53×	1.65×	2.12×	1.91×	1.69×

TABLE VIII

RUNTIME COMPARISON AGAINST BASELINE METHODS ON THREE LARGE DATASETS. AVERAGE TRAINING TIME PER EPOCH [SEC] IS REPORTED. BASELINE REFERS TO GRAPHSAINT FOR *Reddit* AND *ogbn-products*. GAT IS THE BASELINE FOR *ogbn-proteins*.

	<i>Reddit</i>	<i>ogbn-products</i>	<i>ogbn-proteins</i>
Baseline	2.02	170.75	269.70
SUGAR	0.88	77.05	142.7
Speedup	2.30×	2.22×	1.89×

TABLE IX

EVALUATIONS OF SUGAR ON NVIDIA JETSON NANO FOR *Flickr* AND *ogbn-arxiv*. ‘AVG. TIME’ AND ‘MAX MEM’ DENOTE TRAINING TIME PER EPOCH AND PEAK RESIDENT SET SIZE (RSS) MEMORY. WE MEASURE THE TIME, MEMORY AND ENERGY FOR TRAINING 10 EPOCHS. SUGAR IMPROVES AVERAGE TRAINING TIME, MEMORY USAGE AND ENERGY CONSUMPTION PER DEVICE OVER BASELINE GNNs (*i.e.*, GRAPHSAINT AND GCN).

Dataset		Avg. Time [sec]	Max Mem [GB]	Energy [kJ]
<i>Flickr</i>	GraphSAINT-N	22.62	1.05	1.13
	SUGAR	10.50	0.89	0.52
	Improvement	2.15×	1.18×	2.17×
<i>ogbn-arxiv</i>	GCN	28.10	2.24	1.27
	SUGAR	18.39	1.46	0.81
	Improvement	1.53×	1.53×	1.57×

memory usage. In contrast, SUGAR effectively alleviates the issue with 1.52×

 runtime speedup and 1.73× memory reduction.

2) Evaluations on mobile and edge devices:

Following the GPU setting, we proceed to evaluate SUGAR on mobile and edge devices with CPUs.

Training Time. Table VII presents the average training time per epoch of SUGAR compared with

baselines on the *Flickr* and *ogbn-arxiv* datasets. Due to the relative small size of these two datasets, we are able to train GNNs on all five hardware devices, ranging from a Raspberry Pi 3B, to a desktop equipped with high-performance CPUs. We also list the runtime on GPUs in the last column for easy comparison.

From Table VII, we can see that SUGAR demonstrates consistent speedup across all platforms, achieving over 2×

 and 1.5× speedup on the *Flickr* and *ogbn-arxiv* datasets, respectively. In addition, training a GCN on the Raspberry Pi 3B fails due to running out of memory, while SUGAR demonstrates good memory scalability and hence it can be used with such a device with a limited memory budget (*i.e.*, 1GB in this case). This also holds true for the *Reddit* dataset: SUGAR provides a feasible solution for local training on the Jetson Nano (time per epoch is 50.27s), while other baselines can not work due to large memory requirements.

Thus, for the other three datasets, we compare the runtime on Desktop-CPU and report our results in Table VIII. We also observe consistent speedup across all datasets: SUGAR nearly halves the training time in all three cases.

Memory Usage. We compare the memory usage of SUGAR against GNN baselines on a CPU setting. Figure 2 illustrates the resident set size (RSS) memory usage during training on the four datasets: *ogbn-arxiv*, *Reddit*, *ogbn-proteins* and *ogbn-products*. Note that we train a full-batch version of GCN and the batch size of GAT is larger compared with GraphSAGE and GraphSAINT. This accounts for higher fluctuation in the corresponding figure. It is evident that our proposed SUGAR achieves substantial memory reductions compared with baseline GNNs.

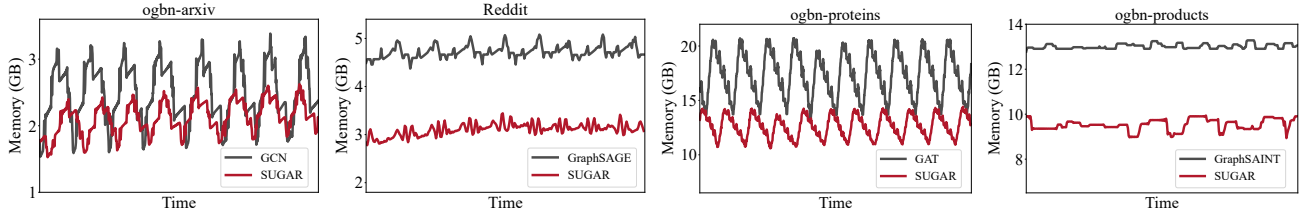


Fig. 2. Memory variation during training GNNs on Desktop-CPU for *ogbn-arxiv*, *Reddit*, *ogbn-proteins* and *ogbn-products*. For SUGAR, we plot the memory variation of the device that consumes most memory.

We emphasize that memory plays a critical role in GNN training. In the context of devices with limited resources, the situation is more severe since the graph dataset is already big and loading the full dataset may not be possible. By adopting subgraph-level training, SUGAR effectively alleviates the problem.

Finally, we present a case study of SUGAR on NVIDIA Jetson Nano in Table IX to demonstrate the applicability of SUGAR to edge devices. Jetson Nano is a popular, cheap and readily available platform (we adopt the model with quad Cortex-A57 CPU and 4GB LPDDR memory) and thus considered as a good fit for our problem scenario. Apart from training time, we measure the peak RSS memory usage for the training process and calculate energy consumption. As shown in Table IX, SUGAR achieves low latency, consumes less memory and is more energy efficient when compared with baseline GNN algorithms. Therefore, it provides an ideal choice to train GNNs on devices with limited memory and battery capacity.

C. Scalability Analysis

1) *Number of partitions*: So far we have demonstrated the great performance of SUGAR with two available devices. A natural follow-up question is, *how does SUGAR perform on more devices, i.e., device number $K > 2$* . Below we provide a scalability analysis of SUGAR based on the number of partitions (*i.e.*, device number K).

We vary the number of available devices K from 2 to 8 and evaluate SUGAR on the *ogbn-arxiv*, *Reddit* and *ogbn-products* datasets. The evaluation is conducted on Desktop-GPU. Runtime speedup, peak GPU memory reduction, validation and test accuracy are presented in Figure 3. With increasing K , we observe a decreased training time and peak memory usage for each local device.

As we can see, while distributing the GNN model to more devices yields computation efficiency, test

accuracy drops a bit. For instance, in the case of 8 devices, the biggest decrease happens in the *ogbn-products* dataset: test accuracy is 76.69% while the baseline accuracy is 79.54%. In the meantime, SUGAR leads to $5.13\times$ speedup, as well as $4.24\times$ memory reduction compared with the baseline. Generally speaking, there exists a tradeoff between training scalability and performance. The underlying reason is that the increase of partition number K leads to more inter-device edges, which corresponds to a larger error in estimating with A_{SG} with A .

We further evaluated SUGAR in a 128-device setting. The results show that the test accuracy drop compared with baseline GNNs is small, *i.e.*, within 5% when scaling up to 128 devices (*e.g.*, accuracy decreases from 72.37% to 67.80% for *ogbn-arxiv*, from 96.39% to 92.32% for *Reddit*, from 50.64% to 46.31% for *Flickr*). At the same time, we note that the memory savings are great (*e.g.*, peak memory usage per device is reduced from 1.60GB to 0.02GB for *ogbn-arxiv*). This shows that SUGAR can work with very small computation and memory requirements at the cost of slightly downgraded performance. Thus, SUGAR provides a feasible solution in extremely resource-limited scenarios while general GNN training methods are not applicable.

2) *Graph Partitioning*: Below we analyze the two graph partitioning strategies introduced in Section III-C: (a) weighted graph partitioning and (b) sub-graph expansion.

Table X compares results achieved with different graph partitioning strategies. As can be seen from the table, compared with random and METIS partitioning, degree-based METIS partitioning (SUGAR) results in a smaller performance degradation. It is worth noting that SUGAR is a flexible algorithm that can incorporate any graph partitioning algorithm. In practice, we adopt degree-based METIS as it is

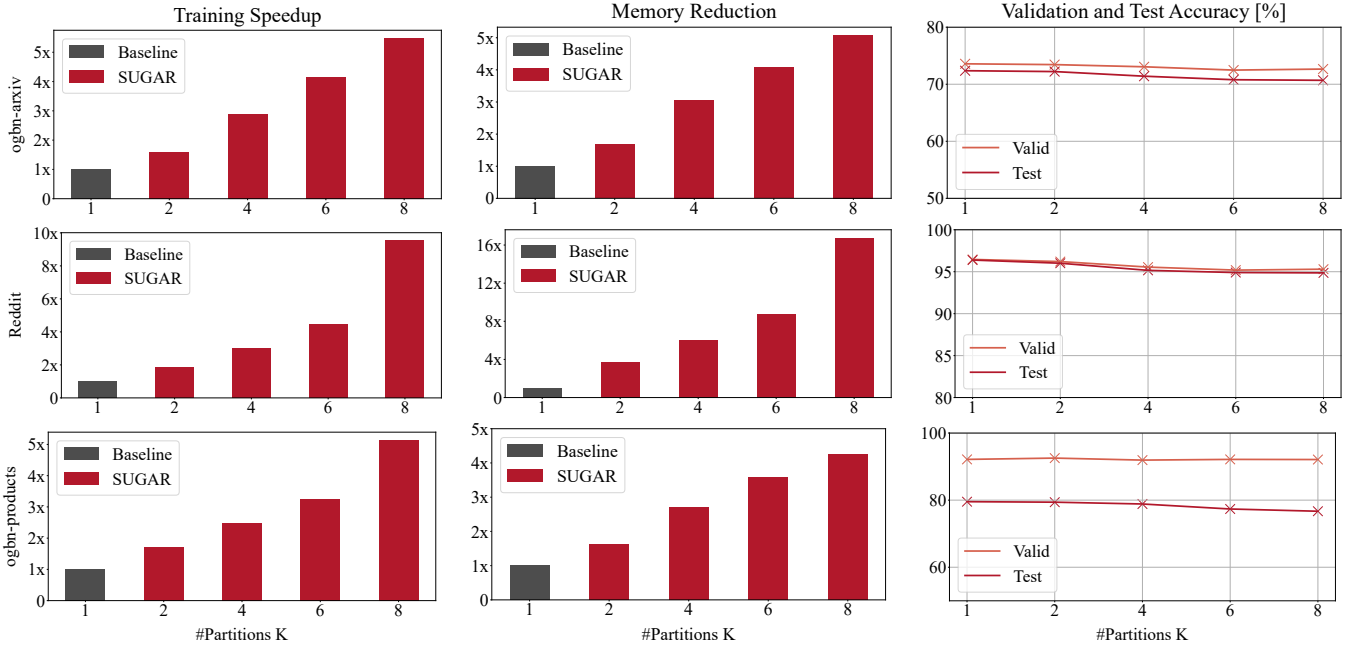


Fig. 3. Scalability analysis on the number of partitions (*i.e.*, the number of available devices K) for SUGAR. $K = 1$ refers to the baseline GNN (*i.e.*, GCN for *ogbn-arxiv*; GraphSAGE for *Reddit*; GraphSAINT for *ogbn-products*). We report the smallest training time speedup and peak GPU memory reduction among K devices (*i.e.*, the worst-case scenario) of SUGAR over the baseline.

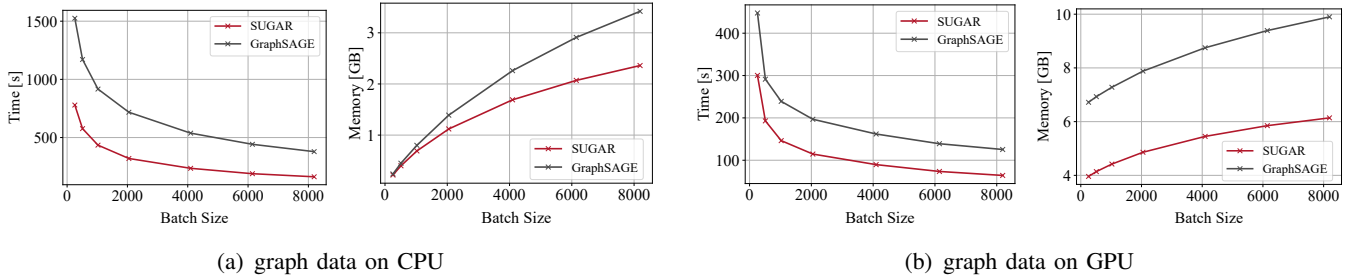


Fig. 4. The training time and peak GPU memory with varying batch sizes of GraphSAGE and SUGAR for the *ogbn-products* data. We investigate two settings: (a) graph data loaded on CPU for memory savings; (b) graph data loaded on GPU for faster execution.

TABLE X
EVALUATIONS OF DIFFERENT GRAPH PARTITIONING STRATEGIES.
WE ADOPT GCN ON THE *ogbn-arxiv* GRAPH, AND BASELINE DENOTES THE RESULTS ACHIEVED WITH THE ORIGINAL GRAPH (*i.e.*, NO GRAPH PARTITIONING).

	Baseline	Random	METIS	SUGAR
Test Acc. [%]	72.37	68.38	70.90	72.22

efficient and easy to implement.

Next, we provide an analysis of the subgraph expansion step. As discussed in Section III, we can expand the node set of each subgraph for better performance, yet the computation and memory costs increase at the same time. We experiment with 4 devices (*i.e.*, $K = 4$) on the *ogbn-arxiv*

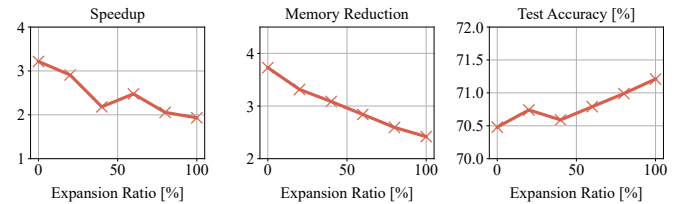


Fig. 5. Results on the *ogbn-arxiv* graph with different graph expansion ratios ($K = 4$). The figure shows the runtime speedup (left), memory reduction (middle) and test accuracy (right) of SUGAR compared with baseline GCN. As shown, subgraph expansion leads to higher accuracy at the cost of increased runtime and memory usage.

graph. To be specific, we expand each subgraph by incorporating $r\%$ 1-hop neighbors that are not originally in the subgraph. Figure 5 provides runtime, memory and accuracy results with r ranging from

0 to 100. A value of 0 represents no subgraph expansion and a value of 100 indicates that all 1-hop neighbors are added to the subgraph. We observe the tradeoff between performance and computational complexity. With increasing r , test performance improves at the cost of increased computation, as well as memory requirements. Thus we consider the subgraph expansion as an optional step that is subject to available resources and users' specific requirements.

3) *Batch Size*: Finally, we study the influence of batch sizes on computational efficiency and memory scalability on SUGAR when compared with mini-batch training algorithms.

For mini-batch training algorithms, when the limited memory of device renders GNN training infeasible, a natural idea is reduce the batch size for memory savings. Here, we analyze the influence of SUGAR and the act of reducing batch sizes on computational efficiency, as well as memory scalability. We conduct experiments on the largest *ogbn-products* graph with GraphSAGE as the baseline. Two settings are considered: (a) graph data loaded on CPU, longer training time and smaller memory consumption is expected; (b) graph data loaded on GPU, the model runs faster, yet requires more GPU memory. Figure 4 provides runtime and memory results with varying batch sizes.

We list two observations below: First, SUGAR mainly improves runtime in setting (a) and achieves greater memory reduction in setting (b). This is related to the mechanism of SUGAR: each local model adopts one subgraph for training instead of the original graph. Thus, data loading time is reduced in setting (a) and putting a subgraph on GPU is more memory efficient in setting (b).

Secondly, SUGAR demonstrates to be a better technique in reducing memory usage than tuning the batch size. While it is generally known that there exists a tradeoff between computation and memory requirements as reducing batch size increases training time, SUGAR is able to improve on both accounts.

V. CONCLUSION

We have proposed SUGAR, an efficient GNN training method that improves training scalability with multiple devices. SUGAR can reduce computation, memory and communication costs during

training through two key contributions: (1) a novel graph partitioning strategy with memory budgets and graph topology taken into consideration; (2) subgraph-level local GNN training. We provided a thorough theoretical analysis of SUGAR and conducted extensive experiments to evaluate SUGAR. Experiments results across multiple hardware platforms demonstrate high training efficiency and memory scalability of SUGAR.

More importantly, SUGAR demonstrates the potential of deploying modern GNN algorithms on resource-limited devices, which opens up discussion in developing resource-efficient GNN approaches that are suitable for IoT deployment. In the future, we plan to extend SUGAR to work with more graph partitioning algorithms and GNN models. One direction is to employ streaming graph partitioning approaches and spatio-temporal GNNs to adapt SUGAR to work with spatio-temporal graphs. Another direction is to address edge-level and graph-level prediction problems with SUGAR.¹

ACKNOWLEDGMENTS

This work was supported in part by the US National Science Foundation (NSF) grants CNS-2007284 and CCF-2107085. The authors would also like to thank Zhengqi Gao (MIT) and Xinyuan Cao (Gatech) for help with theoretical derivations of this work, Guihong Li (UT Austin) and Mengtian Yang (UT Austin) for help with the experiments and the anonymous reviewers for their constructive feedback.

REFERENCES

- [1] Z. Wu, S. Pan, F. Chen, G. Long, C. Zhang, and S. Y. Philip, "A comprehensive survey on graph neural networks," *IEEE Transactions on Neural Networks and Learning Systems*, 2020.
- [2] T. N. Kipf and M. Welling, "Semi-supervised classification with graph convolutional networks," in *International Conference on Learning Representations*, 2017.
- [3] M. Zhang and Y. Chen, "Link prediction based on graph neural networks," *Conference on Neural Information Processing Systems*, 2018.
- [4] W. Fan, Y. Ma, Q. Li, Y. He, E. Zhao, J. Tang, and D. Yin, "Graph neural networks for social recommendation," in *The World Wide Web Conference*, 2019.
- [5] P. Kairouz, H. B. McMahan, B. Avent, A. Bellet, M. Bennis, A. N. Bhagoji, K. Bonawitz, Z. Charles, G. Cormode, R. Cummings *et al.*, "Advances and open problems in federated learning," *arXiv preprint arXiv:1912.04977*, 2019.

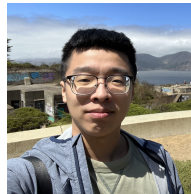
¹See <https://github.com/zihuixue/SUGAR/blob/main/Appendix.pdf> for the Appendix.

- [6] C. Wu, F. Wu, Y. Cao, Y. Huang, and X. Xie, “Fedgcn: Federated graph neural network for privacy-preserving recommendation,” *arXiv preprint arXiv:2102.04925*, 2021.
- [7] H. Bagherinezhad, M. Rastegari, and A. Farhadi, “Lcnn: Lookup-based convolutional neural network,” in *Conference on Computer Vision and Pattern Recognition*, 2017.
- [8] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *Conference on Computer Vision and Pattern Recognition*, 2016.
- [9] W. Hu, M. Fey, M. Zitnik, Y. Dong, H. Ren, B. Liu, M. Catasta, and J. Leskovec, “Open graph benchmark: Datasets for machine learning on graphs,” *Conference on Neural Information Processing Systems*, 2020.
- [10] W. L. Hamilton, R. Ying, and J. Leskovec, “Inductive representation learning on large graphs,” in *Conference on Neural Information Processing Systems*, 2017.
- [11] J. Chen, T. Ma, and C. Xiao, “Fastgcn: fast learning with graph convolutional networks via importance sampling,” *arXiv preprint arXiv:1801.10247*, 2018.
- [12] J. Chen, J. Zhu, and L. Song, “Stochastic training of graph convolutional networks with variance reduction,” *arXiv preprint arXiv:1710.10568*, 2017.
- [13] W.-L. Chiang, X. Liu, S. Si, Y. Li, S. Bengio, and C.-J. Hsieh, “Cluster-gcn: An efficient algorithm for training deep and large graph convolutional networks,” in *International Conference on Knowledge Discovery & Data Mining*, 2019.
- [14] H. Zeng, H. Zhou, A. Srivastava, R. Kannan, and V. Prasanna, “Graphsaint: Graph sampling based inductive learning method,” in *International Conference on Learning Representations*, 2020.
- [15] C. R. Wolfe, J. Yang, A. Chowdhury, C. Dun, A. Bayer, S. Segarra, and A. Kyrillidis, “Gist: Distributed training for large-scale graph convolutional networks,” *arXiv preprint arXiv:2102.10424*, 2021.
- [16] D. Zheng, C. Ma, M. Wang, J. Zhou, Q. Su, X. Song, Q. Gan, Z. Zhang, and G. Karypis, “Distdgl: distributed graph neural network training for billion-scale graphs,” in *Workshop on Irregular Applications: Architectures and Algorithms*, 2020.
- [17] P. Veličković, G. Cucurull, A. Casanova, A. Romero, P. Lio, and Y. Bengio, “Graph attention networks,” *International Conference on Learning Representations*, 2018.
- [18] F. Frasca, E. Rossi, D. Eynard, B. Chamberlain, M. Bronstein, and F. Monti, “Sign: Scalable inception graph neural networks,” *arXiv preprint arXiv:2004.11198*, 2020.
- [19] F. Wu, A. Souza, T. Zhang, C. Fifty, T. Yu, and K. Weinberger, “Simplifying graph convolutional networks,” in *International Conference on Machine Learning*, 2019.
- [20] C. Wan, Y. Li, A. Li, N. S. Kim, and Y. Lin, “Bns-gcn: Efficient full-graph training of graph convolutional networks with partition-parallelism and random boundary node sampling,” *Proceedings of Machine Learning and Systems*, 2022.
- [21] G. Karypis and V. Kumar, “A fast and high quality multilevel scheme for partitioning irregular graphs,” *SIAM Journal on Scientific Computing*, 1998.
- [22] Z. Abbas, V. Kalavri, P. Carbone, and V. Vlassov, “Streaming graph partitioning: an experimental study,” *Proceedings of the VLDB Endowment*, 2018.
- [23] F. Petroni, L. Querzoni, K. Daudjee, S. Kamali, and G. Iacoboni, “Hdfr: Stream-based partitioning for power-law graphs,” in *International Conference on Information and Knowledge Management*, 2015.
- [24] C. Xie, L. Yan, W.-J. Li, and Z. Zhang, “Distributed power-law graph computing: Theoretical and empirical analysis,” *Conference on Neural Information Processing Systems*, 2014.
- [25] D. Luo, W. Cheng, W. Yu, B. Zong, J. Ni, H. Chen, and X. Zhang, “Learning to drop: Robust graph neural network via topological denoising,” in *International Conference on Web Search and Data Mining*, 2021.
- [26] Y. Rong, W. Huang, T. Xu, and J. Huang, “Dropedge: Towards deep graph convolutional networks on node classification,” in *International Conference on Learning Representations*, 2020.
- [27] L. Zhao and L. Akoglu, “Pairnorm: Tackling oversmoothing in gnns,” *arXiv preprint arXiv:1909.12223*, 2019.
- [28] J. Li, T. Zhang, H. Tian, S. Jin, M. Fardad, and R. Zafarani, “Sgcn: A graph sparsifier based on graph convolutional networks,” in *Pacific-Asia Conference on Knowledge Discovery and Data Mining*, 2020.
- [29] T. Chen, Y. Sui, X. Chen, A. Zhang, and Z. Wang, “A unified lottery ticket hypothesis for graph neural networks,” in *International Conference on Machine Learning*, 2021.
- [30] M. Yan, L. Deng, X. Hu, L. Liang, Y. Feng, X. Ye, Z. Zhang, D. Fan, and Y. Xie, “Hygcn: A gcn accelerator with hybrid architecture,” in *International Symposium on High Performance Computer Architecture*, 2020.
- [31] M. Kuczma, *An introduction to the theory of functional equations and inequalities: Cauchy’s equation and Jensen’s inequality*. Springer Science & Business Media, 2009.
- [32] J. Leskovec and A. Krevl, “SNAP Datasets: Stanford large network dataset collection,” <http://snap.stanford.edu/data>, Jun. 2014.
- [33] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga *et al.*, “Pytorch: An imperative style, high-performance deep learning library,” *Conference on Neural Information Processing Systems*, 2019.
- [34] M. Wang, D. Zheng, Z. Ye, Q. Gan, M. Li, X. Song, J. Zhou, C. Ma, L. Yu, Y. Gai *et al.*, “Deep graph library: A graph-centric, highly-performant package for graph neural networks,” *arXiv preprint arXiv:1909.01315*, 2019.

Zihui Xue received the B.S. degree in School of Information Science and Technology, Fudan University, China, in 2020. She is currently pursuing the Ph.D. degree in Electrical and Computer Engineering department at The University of Texas at Austin. Her current research interests include multimodal perception and efficient machine learning.



Yuedong Yang received the B.E. degree in School of Automation Science and Engineering from Xi’an Jiaotong University, China, in 2020. He is currently pursuing the Ph.D. degree in Electrical and Computer Engineering department at The University of Texas at Austin. His current research interests include efficient machine learning algorithms, machine learning systems, and embedded systems.



Radu Marculescu is the Laura Jennings Turner Chair in Engineering and Professor in the Electrical and Computer Engineering department at The University of Texas at Austin. He received his Ph.D. in Electrical Engineering from the University of Southern California in 1998. Radu’s current research focuses on machine learning methods and tools for modeling and optimization of embedded systems, cyber-



physical systems, and social networks. He is a fellow of the ACM and the IEEE.