

---

# Sketch-GNN: Scalable Graph Neural Networks with Sublinear Training Complexity

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Graph Neural Networks (GNNs) are widely applied to graph learning problems  
2 such as node classification. When scaling up the underlying graphs of GNNs to a  
3 larger size, we are forced to either train on the complete graph and keep the full  
4 graph adjacency and node embeddings in memory (which is often infeasible) or  
5 mini-batch sample the graph (which results in exponentially growing computational  
6 complexities with respect to the number of GNN layers). Various sampling-based  
7 and historical-embedding-based methods are proposed to avoid this exponential  
8 growth of complexities. However, none of these solutions eliminates the linear  
9 dependence on graph size. This paper proposes a sketch-based algorithm whose  
10 training time and memory grow sublinearly with respect to graph size by training  
11 GNNs atop a few compact sketches of graph adjacency and node embeddings.  
12 Based on polynomial tensor-sketch (PTS) theory, our framework provides a novel  
13 protocol for sketching non-linear activations and graph convolution matrices in  
14 GNNs, as opposed to existing methods that sketch linear weights or gradients  
15 in neural networks. In addition, we develop a locality sensitive hashing (LSH)  
16 technique that can be trained to improve the quality of sketches. Experiments on  
17 large-graph benchmarks demonstrate the scalability and competitive performance  
18 of our Sketch-GNNs versus their full-size GNN counterparts.

## 19 1 Introduction

20 Graph Neural Networks (GNNs) have achieved the state-of-the-art graph learning in numerous  
21 applications, including classification [26], clustering [3], recommendation systems [42], social  
22 networks [16] and more, through representation learning of target nodes using information aggregated  
23 from neighborhoods in the graph. The manner in which GNNs utilize graph topology, however,  
24 makes it challenging to scale learning to larger graphs or deeper models with desirable computational  
25 and memory efficiency. Full-batch training that stores the Laplacian of the complete graph suffers  
26 from a memory complexity of  $\mathcal{O}(m + ndL + d^2L)$  on an  $n$ -node,  $m$ -edge graph with node features  
27 of dimension  $d$  when employing an  $L$ -layer graph convolutional network (GCN). This linear memory  
28 complexity dependence on  $n$  and the limited memory capacity of GPUs make it difficult to train on  
29 large graphs with millions of nodes or more. As an example, the *MAG240M-LSC* dataset [21] is a  
30 node classification benchmark with over 240 million nodes that takes over 202 GB of GPU memory  
31 when fully loaded.

32 To address the memory constraints, two major lines of research are proposed: (1) Sampling-based  
33 approaches [18, 11, 12, 14, 44] based on the idea of implementing message passing only between  
34 the neighbors within a sampled mini-batch; (2) Historical-embedding based techniques, such as  
35 GNNAutoScale [17] and VQ-GNN [15]), which maintain the expressive power of GNNs on sampled  
36 subgraphs using historical embeddings. However, all of these methods require the number of mini-  
37 batches to be proportional to the size of the graph for fixed memory consumption. In other words,

38 they significantly increase computational time complexity in exchange for memory efficiency when  
 39 scaling up to large graphs. For example, training a 4-layer GCN with just 333K parameters (1.3 MB)  
 40 for 500 epochs on *ogbn-papers100M* can take more than 2 days on a powerful AWS p4d.24x large  
 41 instance [21].

42 We seek to achieve efficient training of GNNs with time and memory complexities sublinear in  
 43 graph size without significant accuracy degradation. Despite the difficulty of this goal, it should be  
 44 achievable given that (1) the number of learnable parameters in GNNs is independent of the graph  
 45 size, and (2) training may not require a traversal of all local neighborhoods on a graph, but rather  
 46 only the most representative ones (thus sublinear in graph size) as some neighborhoods may be  
 47 very similar. In addition, commonly-used GNNs are typically small and shallow with limited model  
 48 capacity and expressive power, indicating that a modest proportion of data may suffice.

49 This paper presents *Sketch-GNN*, a framework for training GNNs with sublinear time and memory  
 50 complexity with respect to graph size. Using the idea of sketching, which maps high-dimensional  
 51 data structures to a lower dimension through entry hashing, we sketch the  $n \times n$  adjacency matrix  
 52 and the  $n \times d$  node feature matrix to a few  $c \times c$  and  $c \times d$  sketches respectively before training,  
 53 where  $c$  is the sketch dimension. While most existing literature focuses on sketching linear weights  
 54 or gradients, we introduce a method for sketching non-linear activation units using polynomial tensor  
 55 sketch theory [19]. This preserves prediction accuracy while avoiding the need to “unsketch” back to  
 56 the original high dimensional graph-node space  $n$ , thereby eliminating the dependence of training  
 57 complexity on the underlying graph size  $n$ . Moreover, we propose to learn and update the sketches in  
 58 an online manner using learnable locality sensitive hashing (LSH) [9]. This reduces the performance  
 59 loss by adaptively enhancing the sketch quality while incurring minor overhead sublinear in graph size.  
 60 In practice, we find that the sketch-ratio  $c/n$  required to maintain “full-graph” model performance  
 61 drops as  $n$  increases; as a result, our Sketch-GNN enjoys sublinear training scalability.

62 *Sketch-GNN* applies sketching techniques to GNNs to achieve training complexity sublinear to the  
 63 *data* size. This is fundamentally different from the few existing works which sketch the weights or  
 64 gradients [29, 13, 25, 28, 36] to reduce the memory footprint of the model and speed up optimization.  
 65 Our approach is flexible to architecture and has the potential to be generalized to other neural networks  
 66 and data types, e.g., CNNs on gigapixel images. To the best of our knowledge, *Sketch-GNN* is the  
 67 first sub-linear complexity training algorithm for GNNs based on LSH and tensor sketching. The  
 68 sublinear efficiency obtained applies not only to GNNs with a fixed convolution matrix, such as  
 69 GCN [26] and GraphSAGE [18], but also to GNNs with learnable convolution matrices, such as  
 70 GAT [38].

71 Experiments on several large graph datasets, such as *ogbn-product* [21] with 2.45M nodes, demon-  
 72 strate that *Sketch-GNNs* can match the performance of the standard model trained on the complete  
 73 graph, while requiring significantly reduced computations and memory for both fixed (GCN, Graph-  
 74 SAGE) and learnable convolution (GAT) models. For instance, SketchGCN on *ogbn-arxiv* [21]  
 75 is 72% and 55% faster than the corresponding full-graph and sampling-based (GraphSAINT [44])  
 76 baselines, while the pre-processing time is just 14% of overall reduction (when running 500 epochs).

## 77 2 Preliminaries

78 **Basic Notations.** Consider a graph with  $n$  nodes and  $m$  edges. Connectivity is given by the adjacency  
 79 matrix  $A \in \{0, 1\}^{n \times n}$  and features on nodes are represented by the matrix  $X \in \mathbb{R}^{n \times d}$ , where  $d$  is the  
 80 number of features. Given a matrix  $C$ , let  $C_{i,j}$ ,  $C_{i,:}$ , and  $C_{:,j}$  denote its  $(i, j)$ -th entry,  $i$ -th row, and  
 81  $j$ -th column, respectively.  $\odot$  denotes the element-wise (Hadamard) product, whereas  $C^{\odot k}$  represents  
 82 the  $k$ -th order element-wise power.  $\|\cdot\|_F$  is the symbol for the Frobenius norm.  $I_n \in \mathbb{R}^{n \times n}$  denotes  
 83 the identity matrix, whereas  $\mathbf{1}_n \in \mathbb{R}^n$  is the vector whose elements are all ones.  $\text{Med}\{\cdot\}$  represents  
 84 the element-wise median over a set of matrices. Superscripts are used to indicate multiple instances  
 85 of the same kind of variable; for instance,  $X^{(l)} \in \mathbb{R}^{n \times d_l}$  are the node representations on layer  $l$ .

86 **Unified Framework of GNNs.** A *Graph Neural Network (GNN)* layer receives the node representa-  
 87 tion of the preceding layer  $X^{(l)} \in \mathbb{R}^{n \times d}$  as input and outputs a new representation  $X^{(l+1)} \in \mathbb{R}^{n \times d}$ ,  
 88 where  $X = X^{(0)} \in \mathbb{R}^{n \times d}$  are the input features. Although GNNs are designed following different  
 89 guiding principles, such as neighborhood aggregation (GraphSAGE), spatial convolution (GCN),  
 90 self-attention (GAT), and Weisfeiler-Lehman (WL) alignment (GIN [43]), the great majority of  
 91 GNNs can be interpreted as performing message passing on node features, followed by feature

92 transformation and an activation function. The update rule of these GNNs can be summarized as [15]

$$X^{(l+1)} = \sigma\left(\sum_q C^{(q)} X^{(l)} W^{(l,q)}\right). \quad (1)$$

93 Where  $C^{(q)} \in \mathbb{R}^{n \times n}$  denotes the  $q$ -th convolution matrix that defines the message passing operator,  
 94  $q \in \mathbb{Z}_+$  is index of convolution,  $\sigma(\cdot)$  is some choice of nonlinear activation function, and  $W^{(l,q)} \in$   
 95  $\mathbb{R}^{d_l \times d_{l+1}}$  denotes the learnable linear weight matrix for the  $l$ -th layer and  $q$ -th filter. GNNs under this  
 96 paradigm differ from each other by their choice of convolution matrices  $C^{(q)}$ , which can be either  
 97 fixed (GCN and GraphSAGE) or learnable (GAT). In Appendix A.1, we re-formulate a number of  
 98 well-known GNNs under this framework. Unless otherwise specified, we assume  $q = 1$  and  $d = d_l$   
 99 for every layer  $l \in [L]$  for notational convenience.

100 **Count Sketch and Tensor Sketch.** (1) *Count sketch* [7, 40] is an efficient dimensionality re-  
 101 duction method that projects an  $n$ -dimensional vector  $\mathbf{u}$  into a smaller  $c$ -dimensional space us-  
 102 ing a random hash table  $h : [n] \rightarrow [c]$  and a binary Rademacher variable  $s : [n] \rightarrow \{\pm 1\}$ ,  
 103 where  $[n] = \{1, \dots, n\}$ . Count sketch is defined as  $\text{CS}(\mathbf{u})_i = \sum_{h(j)=i} s(j)\mathbf{u}_j$ , which is a  
 104 linear transformation of  $\mathbf{u}$ , i.e.,  $\text{CS}(\mathbf{u}) = R\mathbf{u}$ . Here,  $R \in \mathbb{R}^{c \times n}$  denotes the so-called *count*  
 105 *sketch matrix*, which has exactly one non-zero element per column. (2) *Tensor sketch* [31] is  
 106 proposed as a generalization of count sketch to the tensor product of vectors. Given  $\mathbf{z} \in \mathbb{R}^n$   
 107 and an order  $k$ , consider a  $k$  number of i.i.d. hash tables  $h^{(1)}, \dots, h^{(k)} : [n] \rightarrow [c]$  and i.i.d.  
 108 binary Rademacher variables  $s^{(1)}, \dots, s^{(k)} : [n] \rightarrow \{\pm 1\}$ . Tensor sketch also projects vector  
 109  $\mathbf{z} \in \mathbb{R}^n$  into  $\mathbb{R}^c$ , and is defined as  $\text{TS}_k(\mathbf{z})_i = \sum_{h(j_1, \dots, j_k)=i} s^{(1)}(j_1) \dots s^{(k)}(j_k) \mathbf{z}_{j_1} \dots \mathbf{z}_{j_k}$ , where  
 110  $h(j_1, \dots, j_k) = (h^{(1)}(j_1) + \dots + h^{(k)}(j_k)) \bmod c$ . By definition, a tensor sketch of order  $k = 1$   
 111 degenerates to count sketch;  $\text{TS}_1(\cdot) = \text{CS}(\cdot)$ . (3) We define *count sketch of a matrix*  $U \in \mathbb{R}^{d \times n}$  as the  
 112 count sketch of each row vector individually, i.e.,  $\text{CS}(U) \in \mathbb{R}^{d \times c}$  where  $[\text{CS}(U)]_{i,:} = \text{CS}(U_{i,:})$ . The  
 113 *tensor sketch of a matrix* is defined in the same way. Pham and Pagh [31] devise a fast computation  
 114 of tensor sketch of  $U \in \mathbb{R}^{d \times n}$  (sketch dimension  $c$  and order  $k$ ) using count sketches and the Fast  
 115 Fourier Transform (FFT):

$$\text{TS}_k(U) = \text{FFT}^{-1}\left(\bigodot_{p=1}^k \text{FFT}(\text{CS}^{(p)}(U))\right), \quad (2)$$

116 where  $\text{CS}^{(p)}(\cdot)$  is the count sketch with hash function  $h^{(p)}$  and Rademacher variable  $s^{(p)}$ .  $\text{FFT}(\cdot)$   
 117 and  $\text{FFT}^{-1}(\cdot)$  are the FFT and its inverse applied to each row of a matrix.

118 **Locality Sensitive Hashing.** The definition of count sketch and tensor sketch is based on hash  
 119 table(s) that only requires a data independent uniformity, i.e., with high probability the hash-buckets  
 120 are of similar size. In contrast, locality sensitive hashing (LSH) is a hashing scheme that uses  
 121 locality-sensitive hash function  $H : \mathbb{R}^d \rightarrow [c]$  to ensure that nearby vectors are hashed into the  
 122 same bucket (out of  $c$  buckets in total) with high probability while distant ones are not. *SimHash*  
 123 achieves the locality-sensitive property by employing random projections [8]. Given a random matrix  
 124  $P \in \mathbb{R}^{c/2 \times d}$ , SimHash defines a locality-sensitive hash function

$$H(\mathbf{u}) = \arg \max ([P\mathbf{u} \parallel -P\mathbf{u}]), \quad (3)$$

125 where  $[\cdot \parallel \cdot]$  denotes concatenation of two vectors and  $\arg \max$  returns the index of the largest  
 126 element. SimHash is efficient for large batches of vectors [1]. In this paper, we apply a learnable  
 127 version of SimHash that is proposed by Chen et al. [9], in which the projection matrix  $P$  is updated  
 128 using gradient descent; see Section 3.3 for details.

### 129 3 Sketch-GNN Framework via Polynomial Tensor Sketch

130 **Problem and Insights.** We intend to develop a “sketched counterpart” of GNNs, where training  
 131 is based solely on (dimensionality-reduced) compact sketches of the convolution and node feature  
 132 matrices, the sizes of which can be set independently of the graph size  $n$ . In each layer, Sketch-GNN  
 133 receives some sketches of the convolution matrix  $C$  and node representation matrix  $X^{(l)}$  and outputs  
 134 some sketches of the node representations  $X^{(l+1)}$ . As a result, the memory and time complexities  
 135 are inherently independent of  $n$ . The bottleneck of this problem is estimating the nonlinear activated  
 136 product  $\sigma(CX^{(l)}W^{(l)})$ , where  $W^{(l)}$  is the learnable weight of the  $l$ -th layer.

137 Before considering the nonlinear activation, as a first step, we approximate the linear product  
 138  $CX^{(l)}W^{(l)}$ , using dimensionality reduction techniques such as random projections and low-rank  
 139 decompositions. As a direct corollary of the (distributional) Johnson–Lindenstrauss (JL) lemma [24],  
 140 there exists a projection matrix  $R \in \mathbb{R}^{c \times n}$  such that  $CX^{(l)}W^{(l)} \approx (CR^T)(RX^{(l)}W^{(l)})$  [15].

141 Tensor sketch is one of the techniques that can achieve the JL bound [2]; for an error bound, see  
 142 Lemma 1 in Appendix B.

143 Count sketch offers a good estimation of a matrix product,  $CX^{(l)}W^{(l)} \approx CS(C)CS((X^{(l)}W^{(l)})^T)^T$ .  
 144 While tensor sketch can be used to approximate the power of matrix product, i.e.,  $(CX^{(l)}W^{(l)})^{\odot k} \approx$   
 145  $TS_k(C)TS_k((X^{(l)}W^{(l)})^T)^T$ , where  $(\cdot)^{\odot k}$  is the  $k$ -th order element-wise power. If we combine the  
 146 estimators of element-wise powers of  $CX^{(l)}W^{(l)}$ , we can approximate the (element-wise) activation  
 147  $\sigma(\cdot)$  on  $CX^{(l)}W^{(l)}$ . This technique is known as a *polynomial tensor sketch (PTS)* and is discussed  
 148 in [19]. In this paper, we apply PTS to sketch the message passing of GNNs, including the nonlinear  
 149 activations.

### 150 3.1 Sketch-GNN: Approximated Update Rules

151 **Polynomial Tensor Sketch.** Our goal is to approximate the update rule of GNNs (Eq. (1)) in each  
 152 layer. We first expand the element-wise non-linearity  $\sigma$  as a power series, and then approximate the  
 153 powers using count/tensor sketch, i.e.,

$$X^{(l+1)} = \sigma(CX^{(l)}W^{(l)}) \approx \sum_{k=1}^r c_k (CX^{(l)}W^{(l)})^{\odot k} \approx \sum_{k=1}^r c_k TS_k(C) TS_k((X^{(l)}W^{(l)})^T)^T, \quad (4)$$

154 where the  $k = 0$  term always evaluates to zero as  $\sigma(0) = 0$ . In Eq. (4), coefficients  $c_k$  are  
 155 introduced to enable learning or data-driven selection of the weights when combing the terms of  
 156 different order  $k$ . This allows for the approximation of a variety of nonlinear activation functions,  
 157 such as sigmoid and ReLU. The error of this approximation relies on the precise estimation of the  
 158 coefficients  $\{c_k\}_{k=1}^r$ . To identify the coefficients, Han et al. [19] design a coresets-based regression  
 159 algorithm, which requires at least  $O(n)$  additional time and memory. Since the coefficients  $\{c_k\}_{k=1}^r$   
 160 that achieve the best performance for the classification tasks do not necessarily approximate a  
 161 known activation, we propose learning the coefficients  $\{c_k\}_{k=1}^r$  to optimize the classification loss  
 162 directly using gradient descent with simple  $L_2$  regularization. Experiments indicate that the learned  
 163 coefficients can approximate the sigmoid activation with relative errors comparable to those of the  
 164 coresets-based method; see Fig. 1a in Section 5.

165 **Approximated Update Rules.** The remaining step is to approximate the operations of GNNs using  
 166 PTS (Eq. (4)) on sketches of convolution matrix  $C$  and node representation matrix  $X^{(l)}$ . Consider  $r$   
 167 pairwise-independent count sketches  $\{CS^{(k)}(\cdot)\}_{k=1}^r$  with sketch dimension  $c$ , associated with hash  
 168 tables  $h^{(1)}, \dots, h^{(r)}$  and binary Rademacher variables  $s^{(1)}, \dots, s^{(r)}$ , defined prior to training an  
 169  $L$ -layer *Sketch-GNN*. Using these hash tables and Rademacher variables, we may also construct  
 170 tensor sketches  $\{TS_k(\cdot)\}_{k=2}^r$  up to the maximum order  $r$ .

171 In *Sketch-GNN*, sketches of node representations (instead of the  $O(n)$  standard representation) are  
 172 propagated between layers. To get rid of the dependence on  $n$ , we count sketch both sides of Eq. (4)

$$\begin{aligned} S_X^{(l+1,k')} &:= CS^{(k')}((X^{(l+1)})^T) \approx CS^{(k')} \left( \sum_{k=1}^r c_k^{(l)} TS_k((X^{(l)}W^{(l)})^T) TS_k(C)^T \right) \\ &= \sum_{k=1}^r c_k^{(l)} TS_k((X^{(l)}W^{(l)})^T) CS^{(k')} (TS_k(C)^T) \\ &= \sum_{k=1}^r c_k^{(l)} \text{FFT}^{-1} \left( \bigcirc_{p=1}^k \text{FFT}((W^{(l)})^T S_X^{(l,p)}) \right) S_C^{(l,k,k')}, \end{aligned} \quad (5)$$

173 where  $S_X^{(l+1,k')} = CS^{(k')}((X^{(l+1)})^T) \in \mathbb{R}^{d \times c}$  is the transpose of column-wise count sketch of  
 174  $X^{(l+1)}$ , and the superscripts of  $S_X^{(l+1,k')}$  indicate that it is the  $k'$ -th count sketch of  $X^{(l+1)}$  (i.e.,  
 175 sketched by  $CS^{(k')}(\cdot)$ ). In the second line of Eq. (5), we can move the matrix,  $c_k^{(l)} TS_k((X^{(l)}W^{(l)})^T)$ ,  
 176 multiplied on the left to  $TS_k(C)^T$  out of the count sketch function  $CS^{(k')}(\cdot)$ , since the operation of  
 177 row-wise count sketch  $CS^{(k')}(\cdot)$  is equivalent to multiplying the associated count sketch matrix  $R^{(k')}$   
 178 on the right, i.e., for any  $U \in \mathbb{R}^{n \times n}$ ,  $CS^{(k')}(U) = UR^{(k')}$ . In the third line of Eq. (5), we denote the  
 179 “two-sided sketch” of the convolution matrix as  $S_C^{(l,k,k')} := CS^{(k')}(TS_k(C)^T) \in \mathbb{R}^{c \times c}$  and expand  
 180 the tensor sketch  $TS_k((X^{(l)}W^{(l)})^T)$  using the FFT-based formula (Eq. (2)).

181 Eq. (5) is the (recursive) **update rule** of *Sketch-GNN*, which approximates the operation of the  
 182 original GNN and learns the sketches of representations. Looking at the both ends of Eq. (5), we  
 183 obtain a formula that approximates the sketches of  $X^{(l+1)}$  using the sketches of  $X^{(l)}$  and  $C$ , with  
 184 learnable weights  $W^{(l)} \in \mathbb{R}^{d \times d}$  and coefficients  $\{c_k^{(l)} \in \mathbb{R}\}_{k=1}^r$ . The forward-pass and backward-  
 185 propagation between the input sketches  $\{S_X^{(l,k)}\}_{k=1}^r$  and the sketches of the final layer representations  
 186  $\{S^{(L,k)}\}_{k=1}^r$  take  $O(c)$  time and memory (see Section 3.3 for complexity details).

### 187 3.2 Error Bound on Estimated Representation

188 Based on Lemma 1 and the results in [19], we establish an error bound on the estimated final layer  
 189 representation  $\tilde{X}^{(L)}$  for GCN; see Appendix B for the proof and discussions.

190 **Theorem 1.** *For a Sketch-GNN with  $L$  layers, the estimated final layer representation is*  
 191  $\tilde{X}^{(L)} = \text{Med}\{R^{(k)}S_X^{(L,k)} \mid k = 1, \dots, r\}$ , *where the sketches are recursively computed us-*  
 192 *ing Eq. (5). For  $\Gamma^{(l)} = \max\{5\|X^{(l)}W^{(l)}\|_F^2, (2 + 3^r) \sum_i (\sum_j [X^{(l)}W^{(l)}]_{i,j})^r\}$ , it holds that*  
 193  $\mathbf{E}(\|X^{(L)} - \tilde{X}^{(L)}\|_F^2 / \|X^{(L)}\|_F^2) \leq \prod_{l=1}^L (1 + 2/(1 + c\lambda^{(l)2}/nr\Gamma^{(l)})) - 1$ , *where  $\lambda^{(l)} \geq 0$  is the*  
 194 *smallest singular value of the matrix  $Z \in \mathbb{R}^{nd \times r}$  and  $Z_{:,k}$  is the vectorization of  $(CX^{(l)}W^{(l)})^{\odot k}$ .*  
 195 *Moreover, if  $(c\lambda^{(l)2}/nr\Gamma^{(l)}) \gg 1$  holds true for every layer, the relative error is  $O(L(n/c))$ , which*  
 196 *is proportional to the depth of the model, and inversely proportional to the sketch ratio  $(c/n)$ .*

197 **Remarks.** Despite the fact that in Theorem 1 the error bound grows for smaller sketch ratios  $c/n$ ,  
 198 we observe in experiments that the sketch-ratio required for competitive performance decreases as  
 199  $n$  increases; see Section 5. As for the number of independent sketches  $r$ , we know from Lemma 1  
 200 that the dependence of  $r$  on  $n$  is  $r = \Omega(3^{\log_c n})$  which is negligible when  $n$  is not too small; thus, in  
 201 practice  $r = 3$  is used.

202 The theoretical framework may not completely correspond to reality. Experimentally, the coefficients  
 203  $\{\{c_k^{(l)}\}_{k=1}^r\}_{l=1}^L$  with the highest performance do not necessarily approximate a known activation. We  
 204 defer the challenging problem of bounding the error of sketches and coefficients learned by gradients  
 205 to future studies. Although the error bound is in expectation, we do not train over different sketches  
 206 per iteration due to the instability caused by randomness. Instead, we introduce learnable locality  
 207 sensitive hashing (LSH) in the next section to counteract the approximation limitations caused by the  
 208 fixed number of sketches.

### 209 3.3 A Practical Implementation: Learning Sketches using LSH

210 **Motivations of Learnable Sketches.** In Section 3, we apply polynomial tensor sketch (PTS) to  
 211 approximate the operations of GNNs on sketches of the convolution and feature matrices. Nonetheless,  
 212 the pre-computed sketches are fixed during training, resulting in **two major drawbacks**: (1) The  
 213 performance is limited by the quality of the initial sketches. For example, if the randomly-generated  
 214 hash tables  $\{h^{(k)}\}_{k=1}^r$  have unevenly distributed buckets, there will be more hash collisions and  
 215 consequently worse sketch representations. The performance will suffer because only sketches are  
 216 used in training. (2) More importantly, when multiple Sketch-GNN layers are stacked, the input  
 217 representation  $X^{(l)}$  changes during training (starting from the second layer). Fixed hash tables are  
 218 not tailored to the “changing” hidden representations.

219 We seek a method for efficiently constructing high-quality hash tables tailored for each hidden  
 220 embedding. Locality sensitive hashing (LSH) is a suitable tool since it is data-dependent and  
 221 preserves data similarity by hashing similar vectors into the same bucket. This can significantly  
 222 improve the quality of sketches by reducing the errors due to hash collisions.

223 **Combining LSH with Sketching.** At the time of sketching, the hash table  $h^{(k)} : [n] \rightarrow [c]$  is  
 224 replaced with an LSH function  $H^{(k)} : \mathbb{R}^d \rightarrow [c]$ , for any  $k \in [r]$ . Specifically, in the  $l$ -th layer of a  
 225 Sketch-GNN, we hash the  $i$ -th node to the  $H^{(k)}(X_{i,:}^{(l)})$ -th bucket for every  $i \in [n]$ , where  $X_{i,:}^{(l)}$  is the  
 226 embedding vector of node  $i$ . As a result, we define a data-dependent hash table

$$227 h^{(l,k)}(i) = H^{(k)}(X_{i,:}^{(l)}) \quad (6)$$

228 that can be used for computing the sketches of  $S_X^{(l,k)}$  and  $S_C^{(l,k,k')}$ . This LSH-based sketching can  
 229 be directly applied to sketch the fixed convolution matrix and the input feature matrix. If SimHash  
 230 is used, i.e.,  $H^{(k)}(\mathbf{u}) = \arg \max ([P^{(k)}\mathbf{u} \parallel -P^{(k)}\mathbf{u}])$  (Eq. (3)), an additional  $O(nrc(\log c + d))$   
 231 computational overhead is introduced to hash the  $n$  nodes for the  $r$  hash tables during preprocessing;  
 232 see Appendix F more information. SimHash(es) are implemented as simple matrix multiplications  
 233 that are practically very fast.

234 In order to employ LSH-based hash functions customized to each layer to sketch the hidden repre-  
 235 sentations of a Sketch-GNN (i.e.,  $l = 2, \dots, L - 1$ ), we face **two major challenges**: (1) Unless we  
 236 explicitly unsketch in each layer, the estimated hidden representations  $\tilde{X}^{(l)}$  ( $l = 2, \dots, L - 1$ ) cannot  
 237 be accessed and used to compute the hash tables. However, unsketching any hidden representation,  
 i.e.,  $\tilde{X}^{(l)} = \text{Med}\{R^{(k)}S_X^{(l,k)} \mid k = 1, \dots, r\}$ , requires  $O(n)$  memory and time. We need to come

238 up with an efficient algorithm that updates the hash tables without having to unsketch the complete  
 239 representation. (2) It’s unclear how to change the underlying hash table of a sketch across layers  
 240 without unsketching to the  $n$ -dimensional space, even if we know the most up-to-date hash tables  
 241 suited to each layer.

242 The **challenge (2)**, i.e., changing the underlying hash table of across layers, can be solved by  
 243 maintaining a sparse  $c \times c$  matrix  $T^{(l,k)} := R^{(l,k)}(R^{(l+1,k)})^\top$  for each  $k \in [r]$ , which only requires  
 244  $O(cr)$  memory and time overhead; see Appendix C for more information and detailed discussions.  
 245 We focus on **challenge (1)** for the remainder of this section.

246 **Online Learning of Sketches.** To learn a hash table tailored for a hidden layer using LSH without  
 247 unsketching, we develop an efficient algorithm to update the LSH function using only a size- $|B|$   
 248 subset of the length- $n$  unsketched representations, where  $B$  denotes a subset of nodes we select. This  
 249 algorithm, which we term *online learning of sketches*, is made up of two key parts: (*Part 1*) select a  
 250 subset of nodes  $B \subseteq [n]$  to effectively update the hash table, and (*Part 2*) update the LSH function  
 251  $H(\cdot)$  with a triplet loss computed using this subset.

252 (*1*) *Selection of subset  $B$ :* Because model parameters are updated slowly during neural network  
 253 training, the data-dependent LSH hash tables also changes slowly (this behavior was detailed in [9]).  
 254 The amount of updates to the hash table drops very fast along with training, empirically verified  
 255 in Fig. 1b (left) in Section 5. Based on this insight, we only need to update a small fraction of the  
 256 hash table during training. To identify this subset  $B \subseteq [n]$  of nodes, gradient signals can be used.  
 257 Intuitively, a node representation vector hashed into the wrong bucket will be aggregated with distant  
 258 vectors and lead to larger errors and subsequently larger gradient signals. Specifically, we propose  
 259 finding the candidate set  $B$  of nodes by taking the union of the several buckets with the largest  
 260 gradients, i.e.,  $B = \{i \mid h^{(l,k)}(i) = \arg \max_j [S_X^{(l,k)}]_{j,:} \text{ for some } k\}$ . The memory and overhead  
 261 required to update the entries in  $B$  in the hash table is  $O(|B|)$ .

262 (*2*) *Update of LSH function:* In order to update the projection matrix  $P$  that defines a SimHash  
 263  $H^{(k)} : \mathbb{R}^d \rightarrow [c]$  (Eq. (3)), instead of the  $O(n)$  full triplet loss introduced by [9], we consider a  
 264 sampled version of the triplet loss on the candidate set  $B$  with  $O(|B|)$  complexity, namely

$$\mathcal{L}(H, \mathcal{P}_+, \mathcal{P}_-) = \max \left\{ 0, \sum_{(u,v) \in \mathcal{P}_-} \cos(H(u), H(v)) - \sum_{(u,v) \in \mathcal{P}_+} \cos(H(u), H(v)) + \alpha \right\}, \quad (7)$$

265 where  $\mathcal{P}_+ = \{(\tilde{X}_{i,:}, \tilde{X}_{j,:}) \mid i, j \in B, \langle \tilde{X}_{i,:}, \tilde{X}_{j,:} \rangle > t_+\}$  and  $\mathcal{P}_- = \{(\tilde{X}_{i,:}, \tilde{X}_{j,:}) \mid i, j \in$   
 266  $B, \langle \tilde{X}_{i,:}, \tilde{X}_{j,:} \rangle < t_-\}$  are the similar and dissimilar node-pairs in the subset  $B$ ;  $t_+ > t_-$  and  
 267  $\alpha > 0$  are hyper-parameters. This triplet loss  $\mathcal{L}(H, \mathcal{P}_+, \mathcal{P}_-)$  is used to update  $P$  using gradient  
 268 descent, as described in [9], with a  $O(c|B|d + |B|^2)$  overhead. Experimental validation of this LSH  
 269 update mechanism can be found in Fig. 1b in Section 5.

270 **Avoiding  $O(n)$  in Loss Evaluation.** We can estimate the final layer representation using the  $r$   
 271 sketches  $\{S^{(L,k)}\}_{k=1}^r$ , i.e.,  $\tilde{X}^{(L)} = \text{Med}\{R^{(k)}S_X^{(L,k)} \mid k = 1, \dots, r\}$  and compute the losses of all  
 272 nodes for node classification (or some node pairs for link prediction). However, the complexity of  
 273 loss evaluation is  $O(n)$ , proportional to the number of ground-truth labels. In order to avoid  $O(n)$   
 274 complexity completely, rather than un-sketching the node representation for all labeled nodes, we  
 275 employ the locality sensitive hashing (LSH) technique again for loss calculation so that only a subset  
 276 of node losses are evaluated based on a set of hash tables. Specifically, we construct an LSH hash  
 277 table for each class in a node classification problem, which indexes all of the labeled nodes of this  
 278 class and can be utilized to choose the nodes with poor predictions by leveraging the locality property.  
 279 This technique, introduced in [10], is known as sparse forward-pass and back-propagation, and we  
 280 defer the descriptions to Appendix C.

281 **One-time Preprocessing.** If the convolution matrix  $C$  is fixed (GCN, GraphSAGE), the “two-sided  
 282 sketch”  $S_C^{(l,k,k')} = CS^{(k,k')}(\text{TS}_k(C))^\top \in \mathbb{R}^{c \times c}$  is the same in each layer and may be denoted as  
 283  $S_C^{(k,k')}$ . In addition, all of the  $r^2$  sketches of  $C$ , i.e.,  $\{S_C^{(k,k')} \in \mathbb{R}^{c \times c}\}_{k=1}^r\}_{k'=1}^r$  can be computed  
 284 during the preprocessing phase. If the convolution matrix  $C$  is sparse (which is true for most GNNs  
 285 following Eq. (1) on a sparse graph), we can use the sparse matrix representations for the sketches  
 286  $\{S_C^{(k,k')} \in \mathbb{R}^{c \times c}\}_{k=1}^r\}_{k'=1}^r$ , and the total memory taken by the  $r^2$  sketches is  $O(r^2c(m/n))$  where  
 287  $(2m/n)$  is the average node degree (see Appendix F for details). We also need to compute the  $r$  count  
 288 sketches of the input node feature matrix  $X = X^{(0)}$ , i.e.,  $\{S_X^{(0,k)}\}_{k=1}^r$  during preprocessing, which

289 requires  $O(rcd)$  memory in total. In this regard, we have substituted the input data with compact  
290 graph-size independent sketches (i.e.,  $O(c)$  memory). Although the preprocessing time required to  
291 compute these sketches is  $O(n)$ , it is a one-time cost prior to training, and it is widely known that  
292 sketching is practically very fast.

293 **Complexities of Sketch-GCN.** The theoretical complexities of Sketch-GNN is summarized as  
294 follows, where for simplicity we assume bounded maximum node degree, i.e.,  $m = O(n)$ . **(1)**  
295 **Training Complexity:** (1a) *Forward and backward propagation:*  $O(Lcrd(\log(c) + d + m/n)) =$   
296  $O(c)$  time and  $O(Lr(cd + rm/n)) = O(c)$  memory. (1b) *Hash and sketch update:*  $O(Lr(c +$   
297  $|B|d)) = O(c)$  time and memory. **(2) Preprocessing:**  $O(r(rm + n + c)) = O(n)$  time and  
298  $O(rc(d + rm/n)) = O(c)$  memory. **(3) Inference:**  $O(Ld(m + nd)) = O(n)$  time and  $O(m +$   
299  $Ld(n + d)) = O(n)$  memory (the same as a standard GCN). We defer a detailed summary of the  
300 theoretical complexities of Sketch-GNN to Appendix F.

301 We generalize Sketch-GNN to more GNN models in Appendix D and the pseudo-code which outlines  
302 the complete workflow of Sketch-GNN can be find in Appendix E.

## 303 4 Related Work

304 **Scalable methods for GNNs** can be categorized into four classes, all of them still require linear  
305 training complexities. **(A)** On a large sparse graph with  $n$  nodes and  $m$  edges, the “full-graph”  
306 training of a  $L$ -layer GCN with  $d$ -dimensional (hidden) features per layer requires  $O(m + ndL + d^2L)$   
307 memory and  $O(mdL + nd^2L)$  epoch time. **(B)** Sampling-based methods sample mini-batches from  
308 the complete graph following three schemes: (1) node-wisely sample a subset of neighbors in  
309 each layer to reduce the neighborhood size; (2) layer-wisely sample a set of nodes independently  
310 in each layer; (3) subgraph-wisely sample a subgraph directly and simply forward-pass and back-  
311 propagate on that subgraph. **(B.1)** GraphSAGE [18] samples  $r$  neighbors for each node while ignoring  
312 messages from other neighbors.  $O(br^L)$  nodes are sampled in a mini-batch (where  $b$  is the mini-batch  
313 size), and the epoch time is  $O(ndr^L)$ ; therefore, GraphSAGE is impractical for deep GNNs on a  
314 large graph. FastGCN [12] and LADIES [46] are layer-sampling methods that apply importance  
315 sampling to reduce variance. **(B.2)** The subgraph-wise scheme has the best performance and is  
316 most prevalent. Cluster-GCN [14] partitions the graph into many densely connected subgraphs and  
317 samples a subset of subgraphs (with edges between subgraphs added back) for training per iteration.  
318 GraphSAINT [44] samples a set of nodes and uses the induced subgraph for mini-batch training.  
319 Both Cluster-GCN and GraphSAINT require  $O(mdL + nd^2L)$  epoch time, which is the same as  
320 “full-graph” training, although Cluster-GCN also needs  $O(m)$  pre-processing time. **(C)** Apart from  
321 sampling strategies, historical-embedding-based methods propose mitigating sampling errors and  
322 improving performance using some stored embeddings. GNNAutoScale [17] keeps a snapshot of  
323 all embeddings in CPU memory, leading to a large  $O(ndL)$  memory overhead. VQ-GNN [15]  
324 maintains a vector quantized data structure for the historical embeddings, whose size is independent  
325 of  $n$ . **(D)** Linearized GNNs [41, 4, 32] replace the message passing operation in each layer with a  
326 one-time message passing during preprocessing. They are practically efficient, but the theoretical  
327 complexities remain  $O(n)$ . Linearized models usually over-simplify the corresponding GNN and  
328 limit its expressive power.

329 **Towards sublinear GNNs.** Nearly all existing scalable methods focus on mini-batching the large  
330 graph and resolving the memory bottleneck of GNNs, without reducing the epoch training time.  
331 Few recent work focus on graph compression [22, 23] can also achieve sublinear training time by  
332 coarsening (e.g., using [30]) the graph during preprocessing and training GNNs on the coarsened  
333 graph with fewer nodes and edges. Nevertheless, this strategy suffers from **two issues:** **(1)** Although  
334 graph coarsening is a one-time cost, the memory and time overheads are often worse than  $O(n)$   
335 and can be prohibitively large on graphs with over 100K nodes. Even the fastest graph coarsening  
336 algorithm used by [22] takes more than 68 minutes to process the 233K-node *Reddit* graph [44];  
337 see Table 1. The long preprocessing time renders any training speedups meaningless. **(2)** The test  
338 performance of a model trained on the coarsened graph highly depends on the GNN type. Although  
339 the performance of [22] on GCN is good, significant performance degradations are observed on  
340 GraphSAGE and GAT; see Section 5.

341 We defer discussion of more scalable GNN papers and the broad literature of sketching and LHS for  
342 neural networks to Appendix G.

## 5 Experiments

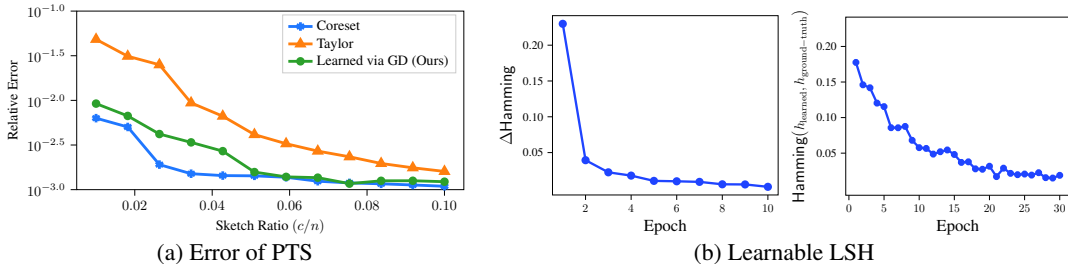


Figure 1: **Figure 1a** Relative errors when applying polynomial tensor sketch (PTS) to the nonlinear unit  $\sigma(CXW)$  following Eq. (4). The dataset used is Cora [33].  $\sigma$  is the sigmoid activation. We set  $r = 5$  and test on a GCN with fixed  $W = I_d \in \mathbb{R}^{d \times d}$ . The coefficients  $\{c_k\}_{k=1}^r$  can be computed by a coresets regression [19] (blue), by a Taylor expansion of  $\sigma(\cdot)$  (orange), or learned from gradient descent proposed by us (green). **Figure 1b** The left plot shows the Hamming distance changes of the hash table in the 2nd layer during the training of a 2-layer *Sketch-GCN*, where the hash table is constructed from the unsketched representation  $\tilde{X}^{(1)}$  using SimHash. The right plot shows the Hamming distances between the hash table learned using our algorithm and the hash table constructed directly from  $\tilde{X}^{(1)}$ .

Table 1: Time and memory efficiencies of Sketch-GNN versus other scalable methods.

| Benchmark                | <i>ogbn-arxiv</i>  |            |              | <i>Reddit</i>      |                  |              |
|--------------------------|--------------------|------------|--------------|--------------------|------------------|--------------|
|                          | Preprocessing Time | Epoch Time | Train Memory | Preprocessing Time | Epoch Time       | Train Memory |
| “Full-Graph” (oracle)    | —                  | 0.49 s     | 983 MB       | —                  | OOM <sup>1</sup> | OOM          |
| GraphSAINT               | —                  | 0.30 s     | 31.4 MB      | —                  | 2.09 s           | 977 MB       |
| VQ-GNN                   | —                  | 0.37 s     | 48.9 MB      | —                  | 2.16 s           | 1281 MB      |
| Coarsening               | 358 s              | 0.20 s     | 22.1 MB      | 4123 s             | 1.04 s           | 530 MB       |
| <b>Sketch-GNN (ours)</b> | 27 s               | 0.13 s     | 38.7 MB      | 141 s              | 0.81 s           | 748 MB       |

<sup>1</sup>“OOM” refers to “out of memory”.

Table 2: Performance of Sketch-GNN in comparison to Graph Coarsening [22] on *ogbn-arxiv*.

| Benchmark                | <i>ogbn-arxiv</i> |               |               |               |               |               |               |               |               |
|--------------------------|-------------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|
|                          | GCN               |               |               | GraphSAGE     |               |               | GAT           |               |               |
| “Full-Graph” (oracle)    | .7174 ± .0029     |               |               | .7149 ± .0027 |               |               | .7233 ± .0045 |               |               |
| Sketch Ratio ( $c/n$ )   | 0.1               | 0.2           | 0.4           | 0.1           | 0.2           | 0.4           | 0.1           | 0.2           | 0.4           |
| Coarsening               | .6508 ± .0091     | .6665 ± .0010 | .6892 ± .0035 | .5264 ± .0251 | .5996 ± .0134 | .6609 ± .0061 | .5177 ± .0028 | .5946 ± .0027 | .6307 ± .0041 |
| <b>Sketch-GNN (ours)</b> | .6913 ± .0154     | .7004 ± .0096 | .7028 ± .0087 | .6929 ± .0194 | .6963 ± .0056 | .7048 ± .0080 | .6967 ± .0067 | .6910 ± .0135 | .7053 ± .0034 |

Table 3: Performance of Sketch-GNN versus SGC [41], GraphSAINT [44], and VQ-GNN [15].

| Benchmark                | <i>ogbn-arxiv</i> |               |               | <i>Reddit</i> |               |               | <i>ogbn-product</i> |               |               |
|--------------------------|-------------------|---------------|---------------|---------------|---------------|---------------|---------------------|---------------|---------------|
|                          | GCN               | GraphSAGE     | GAT           | GCN           | GraphSAGE     | GAT           | GCN                 | GraphSAGE     | GAT           |
| SGC                      | .6944 ± .0005     |               |               | .9464 ± .0011 |               |               | .6683 ± .0029       |               |               |
| “Full-Graph” (oracle)    | .7174 ± .0029     |               |               | OOM           |               |               | OOM                 |               |               |
| GraphSAINT               | .7079 ± .0057     | .6987 ± .0039 | .7117 ± .0032 | .9225 ± .0057 | .9581 ± .0074 | .9431 ± .0067 | .7602 ± .0021       | .7908 ± .0024 | .7971 ± .0042 |
| VQ-GNN                   | .7055 ± .0033     | .7028 ± .0047 | .7043 ± .0034 | .9399 ± .0021 | .9449 ± .0024 | .9438 ± .0059 | .7524 ± .0032       | .7809 ± .0019 | .7823 ± .0049 |
| Sketch Ratio ( $c/n$ )   | 0.4               |               |               | 0.3           |               |               | 0.2                 |               |               |
| <b>Sketch-GNN (ours)</b> | .7028 ± .0087     | .7048 ± .0080 | .7053 ± .0034 | .9280 ± .0034 | .9485 ± .0061 | .9326 ± .0063 | .7659 ± .0086       | .7851 ± .0071 | .7797 ± .0101 |

In this section, we evaluate the proposed *Sketch-GNN* algorithm and compare it with the (oracle) “full-graph” training baseline, a graph-coarsening based method (Coarsening [22]) which has sublinear training time, and other scalable methods including: a sampling-based method (GraphSAINT [44]), a historical-embedding based method (VQ-GNN [15]), and a linearized GNN (SGC [41]). We test on several large graph benchmarks including *ogbn-arxiv* (169K nodes, 1.2M edges), *Reddit* (233K nodes, 11.6M edges), and *ogbn-products* (2.4M nodes, 61.9M edges) from [20, 44]. See Appendix H for the implementation details.

**Proof-of-Concept Experiments: (1) Errors of gradient-learned PTS coefficients:** In Fig. 1a, we train the PTS coefficients to approximate the sigmoid activated  $\sigma(CXW)$  to evaluate its approximation power to the ground-truth activation. The relative errors are comparable to those of the coresets-based method. **(2) Slow-change phenomenon of LSH hash tables:** In Fig. 1b (left), we count the changes of the hash table constructed from an unsketched hidden representation for each epoch, characterized by the Hamming distances between consecutive updates. The changes drop rapidly as training progresses, indicating that apart from the beginning of training, the hash codes of most nodes do not change at each update. **(3) Sampled triplet loss for learnable LSH:** In Fig. 1b (right), we verify the effectiveness of our update mechanism for LSH hash functions, as the learned hash table gradually approaches the “ground truth”, i.e., the hash table constructed from the unsketched hidden representation.



362 **Efficiency of Sketch-GNNs.** For efficiency measures, we are interested in the comparison to  
 363 Coarsening, as both approaches achieve sublinear training time at the cost of some preprocessing  
 364 overheads. We use a 3-layer GCN as the backbone and set the sketch ratios ( $c/n$ , ratio of sketch  
 365 dimension  $c$  to graph size  $n$ ) of both algorithms to  $c/n = 0.1$ , meaning that the coarsened graph  
 366 contains  $n/10$  nodes. We measure their preprocessing time, average epoch training time, and peak  
 367 training memory, as reported in Table 1. Although not rigorously comparable, we also set the mini-  
 368 batch size of GraphSAINT and VQ-GNN to  $b = n/10$ . We report the average epoch training time and  
 369 peak training memory for each method and the "full-graph" training baseline. In addition to Table 1,  
 370 the following are also recorded: (1) Coarsening requires 980 MB to preprocess *ogbn-arxiv*, whereas  
 371 Sketch-GNN only requires 539 MB. (2) Our preprocessing on the largest dataset, *ogbn-product* (2.4M  
 372 nodes), takes only 414s. (3) The wallclock time for the validation accuracy to reach 99% of its best is  
 373  $88 \pm 8$ s for SketchGCN, which is shorter than VQ-GNN’s  $103 \pm 11$ s and GraphSAINT’s  $120 \pm 4$ s.

374 From Table 1 and the aforementioned results, we can draw **four important conclusions**: (1) Sketch-  
 375 GNN achieves the fastest average epoch time. The coarsened graph is typically much denser and  
 376 increases the time required for message passing. (2) Sketch-GNN usually converges faster than  
 377 GraphSAINT and VQ-GNN. (3) Our preprocessing time is significantly less than that of Coarsening.  
 378 Coarsening suffers from an extremely long preprocessing time, rendering the training speed-ups  
 379 meaningless. Moreover, our preprocessing time scales well with graph size and sparsity. (4) We also  
 380 require less preprocessing memory as sketching is linear/multi-linear operation and usually preserves  
 381 sparsity. (5) Sketch-GNN often requires more training memory than Coarsening in order to maintain  
 382 the copies of sketches and additional data structures, although these memory overheads are small.

383 **Performance of Sketch-GNNs.** We first compare the performance of *Sketch-GNN* with Coarsening  
 384 under various sketch ratios to understand how their performance is affected by the memory bottleneck.  
 385 In Table 2, we report the test accuracy of both approaches on *ogbn-arxiv*, with a 3-layer GCN,  
 386 GraphSAGE, or GAT as the backbone and a sketch ratio of 0.1, 0.2, or 0.4. We see there are  
 387 significant performance degradations when applying Coarsening to GraphSAGE and GAT, even under  
 388 sketch ratio 0.4, indicating that Coarsening may be compatible only with specific GNNs (GCN and  
 389 APPNP as explained in [22]). In contrast, the performance drops of Sketch-GNN are always small  
 390 across all architectures, even when the sketch ratio is 0.1. Therefore, our approach generalizes to  
 391 more GNN architectures and consistently outperforms the Coarsening method.

392 We move on to compare Sketch-GNN with linearized GNNs (SGC), sampling-based (GraphSAINT),  
 393 and historical-embedding-based (VQ-GNN) methods. In Table 3, we report the performance of  
 394 SGC, the "full-graph" training (oracle), GraphSAINT and VQ-GNN with mini-batch size 50K  
 395 (their performance is not affected by the choice of mini-batch size if it is not too small), and  
 396 Sketch-GNN with appropriate sketch ratios (0.4 on *ogbn-arxiv*, 0.3 on *Reddit*, and 0.2 on *ogbn-*  
 397 *product*). From Table 3, we confirm that, with an appropriate sketch ratio, the performance of  
 398 *Sketch-GNN* is always close to the "full-graph" oracle and competitive with the other scalable  
 399 approaches. Impressively, the needed sketch ratio  $c/n$  for Sketch-GNN to achieve competitive  
 400 performance reduces as graph size grows. This further illustrates that, as previously indicated,  
 401 the required training complexities (to get acceptable performance) are sublinear to the graph size.

402 **Ablation Studies: (1) Dependence of sketch dimension  $c$  on graph size  $n$ .** Although the theoretical  
 403 approximation error increases under smaller sketch ratio  $c/n$ , we observe competitive experimental  
 404 results with smaller  $c/n$  especially on large graphs. In practice, the sketch-ratio required to maintain  
 405 "full-graph" model performance decreases with  $n$ , as verified in Table 3:  $c/n = 0.4$  is needed on *ogbn-*  
 406 *arxiv* with 169K nodes but  $c/n = 0.2$  is adequate on *ogbn-product* with 2.45M nodes. (2) **Learned**  
 407 **Sketches versus Fixed Sketches.** We find that learned sketches can improve the performance  
 408 of all models and on all datasets. Under sketch-ratio  $c/n = 0.2$ , the Sketch-GCN with learned  
 409 sketches achieves  $0.7004 \pm 0.0096$  accuracy on *ogbn-arxiv* while fixed randomized sketches degrade  
 410 performance to  $0.6649 \pm 0.0106$ .

## 411 6 Conclusion

412 We present *Sketch-GNN*, a sketch-based GNN training framework with sublinear training time and  
 413 memory complexities. Our main contributions are (1) approximating nonlinear operations in GNNs  
 414 using polynomial tensor sketch (PTS) and (2) updating sketches using learnable locality sensitive  
 415 hashing (LSH). Our novel framework has the potential to be applied to other architectures and  
 416 applications where the amount of data makes training even simple models impractical.

417 **References**

- 418 [1] Alexandr Andoni, Piotr Indyk, Thijs Laarhoven, Ilya Razenshteyn, and Ludwig Schmidt.  
419 Practical and optimal lsh for angular distance. *Advances in neural information processing*  
420 *systems*, 28, 2015.
- 421 [2] Haim Avron, Huy Nguyen, and David Woodruff. Subspace embeddings for the polynomial  
422 kernel. *Advances in neural information processing systems*, 27, 2014.
- 423 [3] Filippo Maria Bianchi, Daniele Grattarola, and Cesare Alippi. Spectral clustering with graph  
424 neural networks for graph pooling. In *International Conference on Machine Learning*, pages  
425 874–883. PMLR, 2020.
- 426 [4] Aleksandar Bojchevski, Johannes Klicpera, Bryan Perozzi, Amol Kapoor, Martin Blais, Benedek  
427 Rózemerczki, Michal Lukasik, and Stephan Günnemann. Scaling graph neural networks with  
428 approximate pagerank. In *Proceedings of the 26th ACM SIGKDD International Conference on*  
429 *Knowledge Discovery & Data Mining*, pages 2464–2473, 2020.
- 430 [5] Daniele Calandriello, Alessandro Lazaric, Ioannis Koutis, and Michal Valko. Improved large-  
431 scale graph learning through ridge spectral sparsification. In *International Conference on*  
432 *Machine Learning*, pages 688–697. PMLR, 2018.
- 433 [6] Sarath Chandar, Sungjin Ahn, Hugo Larochelle, Pascal Vincent, Gerald Tesauro, and Yoshua  
434 Bengio. Hierarchical memory networks. *arXiv preprint arXiv:1605.07427*, 2016.
- 435 [7] Moses Charikar, Kevin Chen, and Martin Farach-Colton. Finding frequent items in data streams.  
436 In *International Colloquium on Automata, Languages, and Programming*, pages 693–703.  
437 Springer, 2002.
- 438 [8] Moses S Charikar. Similarity estimation techniques from rounding algorithms. In *Proceedings*  
439 *of the thirty-fourth annual ACM symposium on Theory of computing*, pages 380–388, 2002.
- 440 [9] Beidi Chen, Zichang Liu, Binghui Peng, Zhaozhuo Xu, Jonathan Lingjie Li, Tri Dao, Zhao  
441 Song, Anshumali Shrivastava, and Christopher Re. Mongoose: A learnable lsh framework for  
442 efficient neural network training. In *International Conference on Learning Representations*,  
443 2020.
- 444 [10] Beidi Chen, Tharun Medini, James Farwell, Charlie Tai, Anshumali Shrivastava, et al. Slide: In  
445 defense of smart algorithms over hardware acceleration for large-scale deep learning systems.  
446 *Proceedings of Machine Learning and Systems*, 2:291–306, 2020.
- 447 [11] Jianfei Chen, Jun Zhu, and Le Song. Stochastic training of graph convolutional networks with  
448 variance reduction. In *International Conference on Machine Learning*, pages 942–950. PMLR,  
449 2018.
- 450 [12] Jie Chen, Tengfei Ma, and Cao Xiao. Fastgcn: Fast learning with graph convolutional networks  
451 via importance sampling. In *International Conference on Learning Representations*, 2018.
- 452 [13] Wenlin Chen, James Wilson, Stephen Tyree, Kilian Weinberger, and Yixin Chen. Compressing  
453 neural networks with the hashing trick. In *International conference on machine learning*, pages  
454 2285–2294. PMLR, 2015.
- 455 [14] Wei-Lin Chiang, Xuanqing Liu, Si Si, Yang Li, Samy Bengio, and Cho-Jui Hsieh. Cluster-gcn:  
456 An efficient algorithm for training deep and large graph convolutional networks. In *Proceedings*  
457 *of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*,  
458 pages 257–266, 2019.
- 459 [15] Mucong Ding, Kezhi Kong, Jingling Li, Chen Zhu, John Dickerson, Furong Huang, and Tom  
460 Goldstein. Vq-gnn: A universal framework to scale up graph neural networks using vector  
461 quantization. *Advances in Neural Information Processing Systems*, 34, 2021.
- 462 [16] Wenqi Fan, Yao Ma, Qing Li, Yuan He, Eric Zhao, Jiliang Tang, and Dawei Yin. Graph neural  
463 networks for social recommendation. In *The World Wide Web Conference*, pages 417–426,  
464 2019.

- 465 [17] Matthias Fey, Jan E Lenssen, Frank Weichert, and Jure Leskovec. Gnnautoscale: Scalable and  
466 expressive graph neural networks via historical embeddings. In *International Conference on*  
467 *Machine Learning*, pages 3294–3304. PMLR, 2021.
- 468 [18] William L Hamilton, Rex Ying, and Jure Leskovec. Inductive representation learning on large  
469 graphs. In *Proceedings of the 31st International Conference on Neural Information Processing*  
470 *Systems*, pages 1025–1035, 2017.
- 471 [19] Insu Han, Haim Avron, and Jinwoo Shin. Polynomial tensor sketch for element-wise function of  
472 low-rank matrix. In *International Conference on Machine Learning*, pages 3984–3993. PMLR,  
473 2020.
- 474 [20] Weihua Hu, Matthias Fey, Marinka Zitnik, Yuxiao Dong, Hongyu Ren, Bowen Liu, Michele  
475 Catasta, and Jure Leskovec. Open graph benchmark: Datasets for machine learning on graphs.  
476 *Advances in neural information processing systems*, 33:22118–22133, 2020.
- 477 [21] Weihua Hu, Matthias Fey, Hongyu Ren, Maho Nakata, Yuxiao Dong, and Jure Leskovec.  
478 Ogb-1sc: A large-scale challenge for machine learning on graphs. In *Thirty-fifth Conference on*  
479 *Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)*, 2021.
- 480 [22] Zengfeng Huang, Shengzhong Zhang, Chong Xi, Tang Liu, and Min Zhou. Scaling up graph  
481 neural networks via graph coarsening. In *Proceedings of the 27th ACM SIGKDD Conference*  
482 *on Knowledge Discovery & Data Mining*, pages 675–684, 2021.
- 483 [23] Wei Jin, Lingxiao Zhao, Shichang Zhang, Yozen Liu, Jiliang Tang, and Neil Shah. Graph con-  
484 densation for graph neural networks. In *International Conference on Learning Representations*,  
485 2022.
- 486 [24] William B Johnson and Joram Lindenstrauss. Extensions of lipschitz mappings into a hilbert  
487 space. *Contemporary Mathematics*, 26(189-206):1, 1984.
- 488 [25] Shiva Prasad Kasiviswanathan, Nina Narodytska, and Hongxia Jin. Network approximation  
489 using tensor sketching. In *International Joint Conference on Artificial Intelligence*, pages  
490 2319–2325, 2018.
- 491 [26] Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional  
492 networks. In *International Conference on Learning Representations*, 2017.
- 493 [27] Nikita Kitaev, Lukasz Kaiser, and Anselm Levskaya. Reformer: The efficient transformer. In  
494 *International Conference on Learning Representations*, 2019.
- 495 [28] Yibo Lin, Zhao Song, and Lin F Yang. Towards a theoretical understanding of hashing-based  
496 neural nets. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages  
497 127–137. PMLR, 2019.
- 498 [29] Zirui Liu, Kaixiong Zhou, Fan Yang, Li Li, Rui Chen, and Xia Hu. Exact: Scalable graph neural  
499 networks training via extreme activation compression. In *International Conference on Learning*  
500 *Representations*, 2021.
- 501 [30] Andreas Loukas. Graph reduction with spectral and cut guarantees. *Journal of Machine*  
502 *Learning Research*, 20:1–42, 2019.
- 503 [31] Ninh Pham and Rasmus Pagh. Fast and scalable polynomial kernels via explicit feature maps.  
504 In *Proceedings of the 19th ACM SIGKDD international conference on Knowledge discovery*  
505 *and data mining*, pages 239–247, 2013.
- 506 [32] Emanuele Rossi, Fabrizio Frasca, Ben Chamberlain, Davide Eynard, Michael Bronstein,  
507 and Federico Monti. Sign: Scalable inception graph neural networks. *arXiv preprint*  
508 *arXiv:2004.11198*, 2020.
- 509 [33] Prithviraj Sen, Galileo Namata, Mustafa Bilgic, Lise Getoor, Brian Galligher, and Tina Eliassi-  
510 Rad. Collective classification in network data. *AI magazine*, 29(3):93–93, 2008.

- 511 [34] Yang Shi and Animashree Anandkumar. Higher-order count sketch: Dimensionality reduction  
512 that retains efficient tensor operations. In *2020 Data Compression Conference (DCC)*, pages  
513 394–394. IEEE, 2020.
- 514 [35] Ryan Spring and Anshumali Shrivastava. Scalable and sustainable deep learning via randomized  
515 hashing. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge  
516 Discovery and Data Mining*, pages 445–454, 2017.
- 517 [36] Ryan Spring, Anastasios Kyrillidis, Vijai Mohan, and Anshumali Shrivastava. Compressing  
518 gradient optimizers via count-sketches. In *International Conference on Machine Learning*,  
519 pages 5946–5955. PMLR, 2019.
- 520 [37] Rakshith S Srinivasa, Cao Xiao, Lucas Glass, Justin Romberg, and Jimeng Sun. Fast  
521 graph attention networks using effective resistance based graph sparsification. *arXiv preprint  
522 arXiv:2006.08796*, 2020.
- 523 [38] Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Yoshua  
524 Bengio. Graph attention networks. In *International Conference on Learning Representations*,  
525 2018.
- 526 [39] Yining Wang, Hsiao-Yu Tung, Alexander J Smola, and Anima Anandkumar. Fast and guaranteed  
527 tensor decomposition via sketching. *Advances in neural information processing systems*, 28,  
528 2015.
- 529 [40] Kilian Weinberger, Anirban Dasgupta, John Langford, Alex Smola, and Josh Attenberg. Feature  
530 hashing for large scale multitask learning. In *Proceedings of the 26th annual international  
531 conference on machine learning*, pages 1113–1120, 2009.
- 532 [41] Felix Wu, Amauri Souza, Tianyi Zhang, Christopher Fifty, Tao Yu, and Kilian Weinberger.  
533 Simplifying graph convolutional networks. In *International conference on machine learning*,  
534 pages 6861–6871. PMLR, 2019.
- 535 [42] Shiwen Wu, Fei Sun, Wentao Zhang, Xu Xie, and Bin Cui. Graph neural networks in recom-  
536 mender systems: a survey. *ACM Computing Surveys (CSUR)*, 2020.
- 537 [43] Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural  
538 networks? In *International Conference on Learning Representations*, 2018.
- 539 [44] Hanqing Zeng, Hongkuan Zhou, Ajitesh Srivastava, Rajgopal Kannan, and Viktor Prasanna.  
540 Graphsaint: Graph sampling based inductive learning method. In *International Conference on  
541 Learning Representations*, 2019.
- 542 [45] Cheng Zheng, Bo Zong, Wei Cheng, Dongjin Song, Jingchao Ni, Wenchao Yu, Haifeng Chen,  
543 and Wei Wang. Robust graph representation learning via neural sparsification. In *International  
544 Conference on Machine Learning*, pages 11458–11468. PMLR, 2020.
- 545 [46] Difan Zou, Ziniu Hu, Yewen Wang, Song Jiang, Yizhou Sun, and Quanquan Gu. Layer-  
546 dependent importance sampling for training deep and large graph convolutional networks.  
547 *Advances in Neural Information Processing Systems*, 32:11249–11259, 2019.

## 548 Checklist

- 549 1. For all authors...
- 550 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s  
551 contributions and scope? [Yes] See Section 1.
- 552 (b) Did you describe the limitations of your work? [Yes] Currently, our work has two major  
553 limitations: (1) our theoretical assumptions and results may not perfectly correspond  
554 to the reality; see the theoretical remarks in Section 3, and (2) our implementation is  
555 not fully-optimized with the more advanced libraries; see the efficiency discussions  
556 in Section 5.

- 557 (c) Did you discuss any potential negative societal impacts of your work? [Yes] We see our  
558 work as a theoretical and methodological contribution toward more resource-efficient  
559 graph representation learning. Our methodological advances may enable larger-scale  
560 network analysis for societal good. However, progress in graph embedding learning  
561 may potentially inspire other hostile social network studies, such as monitoring fine-  
562 grained user interactions.
- 563 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
564 them? [Yes]
- 565 2. If you are including theoretical results...
- 566 (a) Did you state the full set of assumptions of all theoretical results? [Yes] See Lemma 1  
567 and Theorem 1.
- 568 (b) Did you include complete proofs of all theoretical results? [Yes] See Appendix B.
- 569 3. If you ran experiments...
- 570 (a) Did you include the code, data, and instructions needed to reproduce the main experi-  
571 mental results (either in the supplemental material or as a URL)? [Yes] See Appendix H.
- 572 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
573 were chosen)? [Yes] See Appendix H.
- 574 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
575 ments multiple times)? [Yes] See Section 5.
- 576 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
577 of GPUs, internal cluster, or cloud provider)? [Yes] See Appendix H.
- 578 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 579 (a) If your work uses existing assets, did you cite the creators? [Yes] See Appendix H.
- 580 (b) Did you mention the license of the assets? [Yes] See Appendix H.
- 581 (c) Did you include any new assets either in the supplemental material or as a URL? [No]
- 582 (d) Did you discuss whether and how consent was obtained from people whose data you're  
583 using/curating? [No]
- 584 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
585 information or offensive content? [No]
- 586 5. If you used crowdsourcing or conducted research with human subjects...
- 587 (a) Did you include the full text of instructions given to participants and screenshots, if  
588 applicable? [N/A]
- 589 (b) Did you describe any potential participant risks, with links to Institutional Review  
590 Board (IRB) approvals, if applicable? [N/A]
- 591 (c) Did you include the estimated hourly wage paid to participants and the total amount  
592 spent on participant compensation? [N/A]