
GARNET: Reduced-Rank Topology Learning for Robust and Scalable Graph Neural Networks

Anonymous Author(s)

Anonymous Affiliation

Anonymous Email

Abstract

Graph neural networks (GNNs) have been increasingly deployed in various applications that involve learning on non-Euclidean data. However, recent studies show that GNNs are vulnerable to graph adversarial attacks. Although there are several defense methods to improve GNN robustness by eliminating adversarial components, they may also impair the underlying clean graph structure that contributes to GNN training. In addition, few of those defense models can scale to large graphs due to their high computational complexity and memory usage. In this paper, we propose GARNET, a scalable spectral method to boost the adversarial robustness of GNN models. GARNET first leverages weighted spectral embedding to construct a base graph, which is not only resistant to adversarial attacks but also contains critical (clean) graph structure for GNN training. Next, GARNET further refines the base graph by pruning additional uncritical edges based on probabilistic graphical model. GARNET has been evaluated on various datasets, including a large graph with millions of nodes. Our extensive experiment results show that GARNET achieves adversarial accuracy improvement and runtime speedup over state-of-the-art GNN (defense) models by up to 10.23% and 14.7 \times , respectively.

1 Introduction

Recent years have witnessed a surge of interest in graph neural networks (GNNs), which incorporate both graph structure and node attributes to produce low-dimensional embedding vectors that maximally preserve graph structural information [1]. GNNs have achieved promising results in various real-world applications, such as recommendation systems [2], self-driving car [3], and chip placements [4]. However, recent studies have shown that adversarial attacks on graph structure accomplished by inserting, deleting, or rewiring edges in an unnoticeable way, can easily fool the GNN models and drastically degrade their accuracy in downstream tasks (e.g., node classification) [5, 6].

In literature, one of the most effective ways to defend GNNs is to purify the graph by removing adversarial graph structures. Entezari et al. [7] observe that adversarial attacks mainly affect high-rank graph properties; thus they propose to first construct a low-rank graph by performing truncated singular value decomposition (TSVD) on the graph adjacency matrix, which can then be exploited for training a robust GNN model. Later, Jin et al. [8] propose Pro-GNN to jointly learn a new graph and a robust GNN model with the low-rank constraints imposed by the graph structure. While prior methods using low-rank approximation largely eliminate adversarial components in the graph spectrum, they involve dense adjacency matrices during GNN training, leading to a much higher time/space complexity and prohibiting their applications in large-scale graph learning tasks.

In addition, due to the high computational cost of TSVD, existing low-rank based methods can only preserve top r singular components (e.g., $r = 50$). Consequently, as shown in Figure 1(a), these methods may lose a wide range of clean graph spectrum that corresponds to important structures of the clean graph in the spatial domain. This is confirmed in Figure 1(c), where the clean accuracy of the TSVD-based method largely increases when preserving more spectral information via increasing the graph rank r . In other words, prior low-rank approximation methods eliminate high-rank adversarial

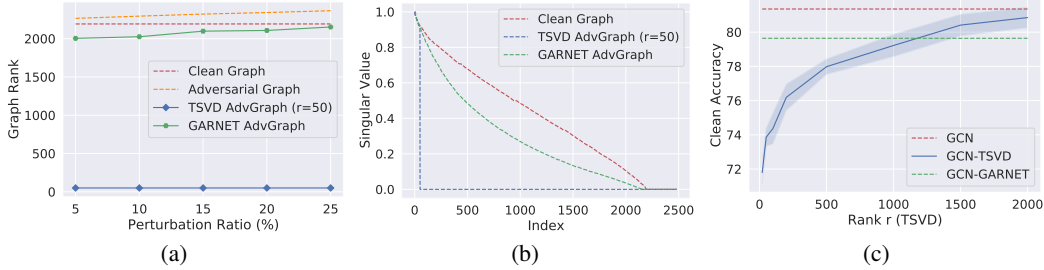


Figure 1: “TSVD AdvGraph” and “GARNET AdvGraph” denote adversarial graphs purified by TSVD and GARNET, respectively. (a) Graph rank comparison on Cora under Metattack with different perturbation ratio. (b) Singular value comparison of different normalized graph adjacency matrices on Cora. (c) Accuracy \pm std. of GCN-TSVD on Cora with different r -rank approximation via TSVD.

41 components at the cost of inevitably impairing the important (clean) graph structure, which degrades
 42 the overall quality of the reconstructed graph and therefore limits the performance of GNN training.

43 In this work, we propose GARNET, a novel spectral approach to learning the underlying clean graph
 44 topology of an adversarial graph via combining spectral embedding with probabilistic graphical
 45 model (PGM), where the learned graph structure encodes the conditional dependence among low-
 46 dimensional node representations (spectral embedding vectors) [9]. More concretely, given an
 47 adversarial graph, GARNET first constructs a base graph topology by leveraging weighted spectral
 48 embeddings that are resistant to adversarial attacks, which is followed by an effective and efficient
 49 graph refinement scheme for pruning noncritical edges in the base graph by exploiting PGM.

50 By recovering the clean graph structure, Figures 1(a) and 1(b) show that the adversarial graph purified
 51 by GARNET largely restores the rank of the underlying clean graph. Thus, GARNET can be viewed
 52 as a reduced-rank topology learning approach that slightly reduces the rank of the input adversarial
 53 graph, which is fundamentally different from the prior low-rank based defense methods (e.g., TSVD
 54 and ProGNN). Moreover, GARNET scales comfortably to large graphs due to its nearly-linear
 55 algorithm complexity, and produces a sparse yet high-quality graph that improves GNN robustness
 56 without involving any dense adjacency matrices during GNN training. As a byproduct, unlike existing
 57 defense methods (e.g., ProGNN) that assume graphs to be homophilic, i.e, adjacent nodes in a graph
 58 tend to have similar attributes [10], GARNET does not have such an assumption and thus can protect
 59 GNNs against adversarial attacks on both homophilic and heterophilic graphs.

60 We evaluate GARNET on both homophilic and heterophilic datasets under strong graph adversarial
 61 attacks such as Nettack [5] and Metattack [6]. Moreover, we further show the nearly-linear scalability
 62 of our approach on the ogbn-products dataset that consists of millions of nodes [11]. Our experimental
 63 results indicate that GARNET largely improves both clean and adversarial accuracy over baselines in
 64 most cases. Our main technical contributions are summarized as follows:

- 65 • To our knowledge, we are the first to exploit spectral graph embedding and probabilistic graphical
 66 model for improving robustness of GNN models, which is achieved by learning a reduced-rank graph
 67 topology for recovering the underlying clean graph structure from the input adversarial graph.
- 68 • By recovering the critical edges that contribute to maximum likelihood estimation in PGM while
 69 ignoring adversarial components, GARNET produces a high-quality graph on which existing GNN
 70 models can be trained to achieve high accuracy. Our experimental results show that GARNET gains
 71 up to 10.23% adversarial accuracy improvement over state-of-the-art defense baselines.
- 72 • Our proposed reduced-rank topology learning method has a nearly-linear complexity in time/space
 73 and produces a sparse graph structure for scalable GNN training. This allows GARNET to run up to
 74 $14.7\times$ faster than prior defense methods on popular data sets such as Cora and Squirrel. In addition,
 75 GARNET scales comfortably to very large graph data sets with millions of nodes, while prior defense
 76 methods run out of memory even on a graph with 20k nodes.

77 2 Background

78 2.1 Undirected Probabilistic Graphical Models

79 Consider an n -dimensional random vector x that follows a multivariate Gaussian distribution $x \sim$
 80 $N(0, \Sigma)$, where $\Sigma = \mathbb{E}[xx^T] \succ 0$ represents the covariance matrix, and $\Theta = \Sigma^{-1}$ represents the

81 precision matrix (inverse covariance matrix). Given a data matrix $X \in R^{n \times d}$ that includes d i.i.d.
 82 (independent and identically distributed) samples $X = [x_1, \dots, x_d]$, where $x_i \sim N(0, \Sigma)$ has an n -
 83 dimensional Gaussian distribution with zero mean, the goal of probabilistic graphical models (PGM)
 84 is to learn a precision matrix Θ that corresponds to an undirected graph structure \mathcal{G} for encoding the
 85 conditional dependence between variables of the observations on columns of X [12, 13]. Specifically,
 86 the classical graphical Lasso method aims at estimating a sparse Θ through maximum likelihood
 87 estimation (MLE) of $f(x)$ leveraging convex optimization [13]. In this work, we focus on one
 88 increasingly popular type of Gaussian graphical models, which is also known as attractive Gaussian
 89 Markov random fields (GMRFs). Attractive GMRFs restrict the precision matrix to be a Laplacian-
 90 like matrix $\Theta = L + \frac{I}{\sigma^2}$, where $L = D - A$ denotes the set of valid graph Laplacian matrices with
 91 D and A representing the diagonal degree matrix and adjacency matrix of the underlying undirected
 92 graph, respectively, I denotes the identity matrix, and σ^2 is a constant denoting prior data variance.
 93 Similar to the graphical Lasso method [13], recent methods for estimating attractive GMRFs leverage
 94 emerging graph signal processing (GSP) techniques to solve the following convex problem [9, 14–17]:
 95

$$\max_{\Theta} \log \det \Theta - \frac{1}{d} \text{tr}(X X^T \Theta) - \alpha \|\Theta\|_1 \quad (1)$$

96 where $\det(\cdot)$ and $\text{tr}(\cdot)$ denote the determinant and trace operators, respectively, α is a hyperparameter
 97 to control the regularization term. The first two terms together can be interpreted as log-likelihood
 98 under a GMRF. The last ℓ_1 regularization term is to enforce Θ (and the corresponding graph) to
 99 be sparse. If X is non-Gaussian, Equation 1 can be regarded as Laplacian estimation based on
 100 minimizing the Bregman divergence between positive definite matrices induced by the function
 101 $\Theta \mapsto -\log \det(\Theta)$ [18].

102 2.2 Graph Adversarial Attacks

103 Most existing graph adversarial attacks aim at degrading the accuracy of GNN models by insert-
 104 ing/deleting edges in an unnoticeable way (e.g., maintaining node degree distribution) [19]. The
 105 most popular graph adversarial attacks fall into the following two categories: (1) targeted attack,
 106 (2) non-targeted attack. The targeted attacks attempt to mislead a GNN model to produce a wrong
 107 prediction on a target sample (e.g., node), while the non-targeted attacks strive to degrade the overall
 108 accuracy of a GNN model for the whole graph data set. Dai et al. [20] first formulate the targeted
 109 attack as a combinatorial optimization problem and leverages reinforcement learning to insert/delete
 110 edges such that the target node is misclassified. Zügner et al. [5] propose another targeted attack
 111 called Nettack, which produces an adversarial graph by maximizing the training loss of GNNs.
 112 Zügner and Günnemann [6] further introduce Metattack, a non-targeted attack that treats the graph
 113 as a hyperparameter and uses meta-gradients to perturb the graph structure. It is worth noting that
 114 graph adversarial attacks have two different settings: poison (perturb a graph prior to GNN training)
 115 and evasion (perturb a graph after GNN training). As shown by Zhu et al. [21], the poison setting is
 116 typically more challenging to defend, as it changes the graph structure that fools GNN training. Thus,
 117 we aim to improve model robustness against attacks under the poison setting.

118 2.3 Graph Adversarial Defenses

119 To defend GNN against adversarial attacks, Entezari et al. [7] first observe that Nettack, a strong
 120 targeted attack, only changes the high-rank information of the adjacency matrix. Thus, they propose
 121 to construct a low-rank graph by performing truncated SVD to undermine the effects of adversarial
 122 attacks. Later, Jin et al. [8] propose Pro-GNN that adopts a similar idea yet encourages nodes with
 123 similar attributes to be connected when jointly learning the low-rank graph and GNN model. Although
 124 those low-rank approximation based methods achieve state-of-the-art results on several datasets,
 125 they produce dense adjacency matrices that correspond to complete graphs, which would limit their
 126 applications for large graphs. Moreover, they only preserve a small region of the graph spectrum
 127 and thus may lose too much important information corresponding to the clean graph structure in the
 128 spatial domain, which limits the performance of GNN training. Recently, [22] exploit Laplacian
 129 eigenpairs to guide GNN training, which produces a robust model with quadratic time complexity
 130 and is thus not scalable to large graphs. In addition to the aforementioned spectral-based defense
 131 methods, GCNJaccard [23] and RS-GNN [24] purify the adversarial graph by connecting nodes with
 132 similar attributes or same labels. However, those defense methods explicitly (or implicitly) assume
 133 the underlying graph to be homophilic, which results in rather poor performance when defending
 134 GNN models on heterophilic graphs. In contrast to the prior arts, GARNET achieves highly robust
 135 yet scalable performance on both homophilic and heterophilic graphs under adversarial attacks by
 136 leveraging a novel graph purification scheme based on spectral embedding and graphical model.

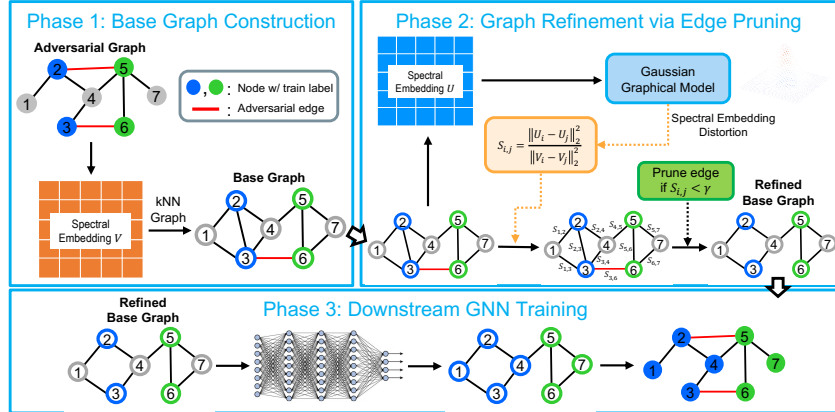
137 **3 The GARNET Approach**

Figure 2: An overview of the three major phases of GARNET.

138 Recently, Entezari et al. [7] and Jin et al. [8] have shown that the well-known graph adversarial
 139 attacks (e.g., Nettack and Metattack) are essentially high-rank attacks, which increase graph rank by
 140 enlarging the smallest singular values of adjacency matrix when perturbing the graph structure, while
 141 rest of the graph spectrum remains almost the same. Consequently, a natural way for improving GNN
 142 robustness is to find the low-rank approximation of the adversarial adjacency matrix.

143 **Low-rank topology learning (prior work).** Given an adversarial adjacency matrix $A_{adv} \in R^{n \times n}$,
 144 Entezari et al. [7] propose to reconstruct a low-rank approximated adjacency matrix via performing
 145 TSVD: $\hat{A} = U\Sigma V^T$, where $\Sigma \in R^{r \times r}$ is a diagonal matrix consisting of r largest singular values
 146 of A_{adv} . $U \in R^{n \times r}$ and $V \in R^{n \times r}$ contain the corresponding left and right singular vectors,
 147 respectively. As the largest singular values are hardly affected by graph adversarial attacks, the
 148 reconstructed low-rank adjacency matrix \hat{A} is resistant to adversarial attacks.

149 However, due to the high computational cost of TSVD, \hat{A} is typically computed by only using top r
 150 largest singular values and their corresponding singular vectors, where r is a relatively small number
 151 (e.g., $r = 50$). Consequently, the rank of \hat{A} is only $r = 50$, which is two orders of magnitude
 152 smaller than the rank of the clean graph, as shown in Figure 1(a). Since these low-rank methods are
 153 overly aggressive in reducing the graph rank, \hat{A} may lose too much important spectral information
 154 corresponding to the clean graph structure. As shown in Figure 1(c), the clean accuracy of the
 155 TSVD-based method is largely improved by increasing the graph rank r , which indicates the low-rank
 156 graph obtained with a small r loses the key graph structure contributing to GNN training. Note that
 157 the adversarial and clean graphs share most of the graph structure, as adversarial attacks perturb the
 158 clean graph in an unnoticeable way. Consequently, losing those important clean graph structures will
 159 also limit the performance of GNN on the adversarial graph.

160 **Reduced-rank topology learning (this work).** Given the adversarial graph \mathcal{G}_{adv} and its adjacency
 161 matrix A_{adv} , our goal is to learn a reduced-rank graph, which slightly reduces the rank of \mathcal{G}_{adv}
 162 to mitigate the effects of adversarial attacks, while retaining most of the important graph spectrum
 163 corresponding to the clean graph structure. As adversarial attacks mainly affect the least dominant
 164 singular components of A_{adv} [7], one straightforward way for constructing such a reduced-rank graph
 165 is to utilize all the singular components except those least dominant ones via TSVD. Nonetheless,
 166 computing such a large number of singular components is computationally expensive [25], and is
 167 thus not scalable to large graphs.

168 To learn the reduced-rank graph in a scalable way, in this work, we leverage only the top few (e.g.,
 169 50) dominant singular components of A_{adv} to restore its important graph spectrum, via recovering
 170 the corresponding clean graph structure with the aid of PGM. Figure 2 gives an overview of our
 171 proposed approach, GARNET, which consists of three major phases. The first phase constructs a
 172 base graph by exploiting spectral embedding and a scalable nearest-neighbor graph algorithm. The
 173 second phase further refines the base graph by pruning noncritical edges based on PGM. The last
 174 phase trains existing GNN models on the refined base graph to improve their robustness. Next, we
 175 will first describe our notion of clean graph recovery via PGM as well as the scalability issue of prior

176 PGM-based work in Section 3.1, which motivates us to develop scalable GARNET kernels described
 177 in Sections 3.2 and 3.3. We further provide the overall complexity of GARNET in Section 3.4.

178 3.1 Graph Recovery via Graphical Model

179 A general philosophy behind PGM is that there exists an underlying graph G , whose structure
 180 determines the joint probability distribution of the observations on the data entities, i.e., columns of a
 181 data matrix $X \in R^{n \times d}$, where n is the number of data points, d the dimension per data point. To
 182 recover the underlying graph structure from the data matrix X , one common way is to leverage MLE
 183 by solving Equation 1 in Section 2.1. As the top few dominant singular components of the adjacency
 184 matrix capture the corresponding graph structure, we can naturally construct the data matrix X based
 185 on those dominant singular components, and then adopt PGM to recover an underlying graph via
 186 MLE. To this end, we define a weighted spectral embedding matrix as follows:

187 **Definition 3.1.** Given the top r smallest eigenvalues $\lambda_1, \lambda_2, \dots, \lambda_r$ and their corresponding eigen-
 188 vectors v_1, v_2, \dots, v_r of normalized graph Laplacian matrix $L_{norm} = I - D^{-\frac{1}{2}}AD^{-\frac{1}{2}}$, where
 189 I and A are the identity matrix and graph adjacency matrix, respectively, and D is a diag-
 190 onal matrix of node degrees, the **weighted spectral embedding matrix** is defined as $V \stackrel{\text{def}}{=} \left[\sqrt{|1 - \lambda_1|}v_1, \dots, \sqrt{|1 - \lambda_r|}v_r \right]$, whose i -th row $V_{i,:}$ is the **weighted spectral embedding** of the
 191 corresponding i -th node in the graph.

193 **Proposition 3.2.** Given a normalized graph adjacency matrix $A_{norm} = D^{-\frac{1}{2}}AD^{-\frac{1}{2}}$ and weighted
 194 spectral embedding matrix V of an undirected graph, let \hat{A} be the rank- r approximation of A_{norm}
 195 via TSVD. If the top r dominant eigenvalues of A_{norm} are non-negative, then we have $\hat{A} = VV^T$.

196 Our proof for Proposition 3.2 is available in Appendix A. Proposition 3.2 shows the connection
 197 between weighted spectral embedding and the low-rank adjacency matrix \hat{A} obtained by TSVD.
 198 Specifically, the weighted spectral embedding matrix V can be viewed as an eigensubspace matrix
 199 consisting of a few dominant singular components of the corresponding adjacency matrix. Thus,
 200 we can use V to recover the underlying clean graph via PGM. However, obtaining V requires the
 201 knowledge of the clean graph structure, which seems to create a chicken and egg problem.

202 Fortunately, since the dominant singular components are hardly affected by adversarial attacks [7],
 203 the weighted spectral embedding is therefore also resistant to adversarial attacks, indicating that the
 204 underlying clean graph \mathcal{G}_{clean} and its corresponding adversarial graph \mathcal{G}_{adv} share almost the same
 205 weighted spectral embeddings. As a result, we can exploit the weighted spectral embedding matrix V
 206 of \mathcal{G}_{adv} to represent that of \mathcal{G}_{clean} . By replacing the data matrix X with V in Equation 1, we have
 207 the following objective function:

$$\max_{\Theta} : F = \log \det \Theta - \frac{1}{r} \text{tr}(VV^T \Theta) - \alpha \|\Theta\|_1 \quad (2)$$

208 [More discussions on Equation 2 are available in Appendix Q.](#) By finding the optimizer Θ^* , we can
 209 recover the underlying graph that maximizes the likelihood given the observation on the weighted
 210 spectral embedding V . However, solving Equation 2 requires at least $O(n^2)$ time/space complexity
 211 per iteration with the most efficient algorithms, which thus cannot scale to large graphs [13, 26, 27].

212 As Θ is constrained to be a Laplacian-like matrix, finding the optimizer Θ^* in Equation 2 is equivalent
 213 to searching for critical edges from a complete graph, which would involve all possible (i.e., $O(n^2)$)
 214 edges. Here we say an edge is critical (noncritical) if including it to the graph significantly increases
 215 (decreases) F in Equation 2. Hence we can recover the underlying graph by pruning noncritical edges
 216 from the complete graph. However, storing a complete graph is still expensive. To have a near-linear
 217 algorithm for clean graph recovery, instead of searching in the complete graph, we limit our search
 218 within an initial base graph \mathcal{G}_{base} that is much sparser but containing sufficient information for
 219 identifying the candidate edges critical to recover the clean graph. Subsequently, the final graphical
 220 model (graph Laplacian) can be obtained by further pruning noncritical edges from \mathcal{G}_{base} .

221 3.2 Base Graph Construction

222 During the first phase of GARNET (shown in Figure 2), our goal is to build a base graph \mathcal{G}_{base} , which
 223 greatly reduces the search space by not constructing a complete graph while preserving the critical
 224 candidate edges that are key to clean graph recovery. To this end, we give the following theorem:

225 **Theorem 3.3.** Given a graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ and its normalized Laplacian matrix $L_{\mathcal{G}}$, let V_i denote
 226 the weighted spectral embedding of node i by using top r eigenpairs of $L_{\mathcal{G}}$. Suppose a relatively

227 *small r is picked such that $\lambda_r \leq 1$, where λ_r is the r -th smallest eigenvalue of L_G , then we have*
 228 $\sum_{(i,j) \in \mathcal{E}} \|V_i - V_j\|_2^2 \leq 0.25r$.

229 Our proof for Theorem 3.3 is available in Appendix B. Note that r is a small constant, which is
 230 independent of the graph size. Thus, Theorem 3.3 indicates that, if an edge connects nodes i and j
 231 in the clean graph, then the Euclidean distance between the weighted spectral embeddings of these
 232 two nodes will be small, which motivates us to build a k -nearest neighbor (kNN) graph as \mathcal{G}_{base} to
 233 incorporate those clean edges. [The connection between kNN and TSVD is provided in Appendix R.](#)

234 Concretely, we first obtain the weighted spectral embedding matrix V of the input adversarial graph
 235 \mathcal{G}_{adv} to represent that of the underlying clean graph \mathcal{G}_{clean} , as V consists of dominant singular
 236 components that are shared by \mathcal{G}_{adv} and \mathcal{G}_{clean} [7]. We then leverage V to construct a kNN graph,
 237 where each node is connected to its k most similar nodes based on the Euclidean distance between
 238 their spectral embeddings. In this work, we exploit an approximate kNN algorithm for constructing
 239 the graph, which has $O(|\mathcal{V}| \log |\mathcal{V}|)$ complexity and thus can scale to very large graphs [28]. By
 240 choosing a proper k (e.g., $k = 50$), \mathcal{G}_{base} is likely to cover edges in the underlying clean graph. Thus,
 241 \mathcal{G}_{base} can serve as a reasonable search space for identifying critical edges in the next step.

242 3.3 Graph Refinement via Edge Pruning

243 For the second phase of GARNET shown in Figure 2, we refine G_{base} by aggressively pruning
 244 noncritical edges from G_{base} , such that the refined graph only preserves the most important edges
 245 that contribute most to the log-likelihood F in Equation 2.

246 To identify critical (noncritical) edges that can most effectively increase (decrease) F , we exploit
 247 the update of Θ based on gradient ascent: $\Theta \leftarrow \Theta + \eta \frac{\partial F}{\partial \Theta}$, where η is the step size. As mentioned in
 248 Section 2.1, Θ is constrained to be $L + \frac{I}{\sigma^2}$, which means the off-diagonal elements in Θ correspond
 249 to negative of edge weights in the underlying graph, i.e., $\Theta_{i,j} = -w_{i,j}$. Thus, the update of $\Theta_{i,j}$
 250 during gradient ascent can be viewed as:

$$\Theta_{i,j} \leftarrow \Theta_{i,j} + \eta \left(\frac{\partial F}{\partial \Theta} \right)_{i,j} = \Theta_{i,j} - \eta \frac{\partial F}{\partial w_{i,j}} \quad (3)$$

251 Equation 3 means that, if $\frac{\partial F}{\partial w_{i,j}}$ is large and positive, $\Theta_{i,j}$ will become more negative, which
 252 corresponds to increasing the edge weight in the underlying graph. Similarly, if $\frac{\partial F}{\partial w_{i,j}}$ is small and
 253 negative, $\Theta_{i,j}$ will be less negative, corresponding to decreasing the edge weight. In other words,
 254 the edge weight $w_{i,j}$ with a large (small) $\frac{\partial F}{\partial w_{i,j}}$ should be increased (decreased) to maximize the
 255 log-likelihood F , meaning the corresponding edge is critical (noncritical). Thus, we can identify the
 256 critical edges once we know $\frac{\partial F}{\partial w_{i,j}}$. By setting $\alpha = 0$ in Equation 2 (as GARNET naturally produces
 257 a sparse graph) and taking the partial derivative with respect to an edge weight $w_{i,j}$, we have:

$$\frac{\partial F}{\partial w_{i,j}} = \sum_{k=1}^n \frac{1}{\lambda_k + 1/\sigma^2} \frac{\partial \lambda_k}{\partial w_{i,j}} - \frac{\|V^T e_{i,j}\|_2^2}{r} \quad (4)$$

258 where $\lambda_k, \forall k = 1, 2, \dots, n$ are the Laplacian eigenvalues of \mathcal{G}_{base} (the initial graph for edge pruning),
 259 $e_{i,j} = e_i - e_j$, and e_i denotes the vector with all zero entries except for the i -th entry being 1.

260 **Theorem 3.4** (Feng [17]). *Let λ_k and u_k be the k -th eigenvalue and the corresponding eigenvector*
 261 *of the Laplacian matrix, respectively. The spectral perturbation $\delta \lambda_k$ due to the increase of an edge*
 262 *weight $w_{i,j}$ can be estimated by $\delta \lambda_k = \delta w_{i,j} (u_k^T e_{i,j})^2$.*

263 The proof for Theorem 3.4 is available in Feng [17]. According to Theorem 3.4 and Equation 4, we
 264 can estimate $\frac{\partial F}{\partial w_{i,j}} \approx \|U^T e_{i,j}\|_2^2 - \frac{1}{r} \|V^T e_{i,j}\|_2^2$, where $U = \left[\frac{u_1}{\sqrt{\lambda_1 + 1/\sigma^2}}, \dots, \frac{u_r}{\sqrt{\lambda_r + 1/\sigma^2}} \right]$, λ_i is the

265 i -th smallest Laplacian eigenvalue of \mathcal{G}_{base} , and u_i is the corresponding eigenvector. Consequently,
 266 an edge (i, j) is critical if $\|U^T e_{i,j}\|_2^2 \gg \frac{1}{r} \|V^T e_{i,j}\|_2^2$. As V and U are the spectral embeddings
 267 on the input adversarial graph and the base graph, respectively, we define the **spectral embedding**

268 **distortion** $s_{i,j} = \frac{\|U^T e_{i,j}\|_2^2}{\|V^T e_{i,j}\|_2^2}$ to measure the edge importance. Consequently, we prune edges in
 269 the base graph \mathcal{G}_{base} that have small spectral embedding distortion, i.e., $s_{i,j} < \gamma$, where γ is a
 270 hyperparameter to control the sparsity of the refined graph. Hence, the refined base graph \mathcal{G}'_{base}
 271 largely recovers the underlying clean graph structure from the input adversarial graph. Since \mathcal{G}'_{base}
 272 is constructed by only leveraging the top few dominant singular components of \mathcal{G}_{adv} , it ignores the
 273 high-rank adversarial components and thus robust to adversarial attacks. As a result, we can train a
 274 given GNN model on \mathcal{G}'_{base} to improve its robustness, which is the last phase of GARNET.

275 3.4 Complexity of GARNET

276 The first phase of GARNET requires $O(r|\mathcal{E}|)$ time for computing top r Laplacian eigenpairs [25],
 277 and $O(|\mathcal{V}| \log |\mathcal{V}|)$ time for kNN graph construction [28]. The second phase involves $O(rk|\mathcal{V}|)$ time
 278 for computing spectral embeddings and edge pruning on the kNN graph. Thus, the overall time
 279 complexity for graph purification is $O(r(|\mathcal{E}| + k|\mathcal{V}|) + |\mathcal{V}| \log |\mathcal{V}|)$, where $|\mathcal{V}|$ ($|\mathcal{E}|$) denotes the
 280 number of nodes (edges) in the adversarial graph, and k is the averaged node degree in the kNN
 281 graph. Our systematic approach of choosing r and the space complexity analysis are in Appendix F.

282 4 Experiments

283 We have conducted comparative evaluation of GARNET against state-of-the-art defense GNN models
 284 under targeted attack (Nettack) [5] and non-targeted attack (Metattack) [6] on both homophilic
 285 and heterophilic datasets. Besides, we also evaluate GARNET robustness against adaptive attacks.
 286 In addition, we further show the scalability of GARNET by comparing its run time with prior
 287 defense methods and evaluating GARNET on ogbn-products, which consists of more than 2 million
 288 nodes [11]. Finally, we conduct ablation studies to understand the effectiveness of GARNET kernels.

289 **Experimental Setup.** The details of datasets used in our experiments are available in Appendix C.
 290 We choose as baselines two state-of-the-art defense methods based on graph purification: TSVD [7]
 291 and Pro-GNN [8]. Besides, we evaluate training based defense methods GCN-LFR [22] and GN-
 292 NGuard [29] on homophilic and heterophilic graphs, respectively. Moreover, we use GCN [30]
 293 and GPRGNN [31] as the backbone GNN models for defense on homophilic datasets (i.e., Cora
 294 and Pubmed). As GCN performs poorly on heterophilic datasets [10, 32], we choose GPRGNN
 295 as the backbone model (as well as the surrogate model for attacking) on Chameleon and Squirrel
 296 datasets. Due to the space limit, we provide defense results with H2GCN [10] as the backbone model
 297 in Appendix J. For all baselines, we tune their hyperparameters against adversarial attacks with a
 298 small perturbation, and keep the same hyperparameters for larger adversarial perturbations. Detailed
 299 hyperparameter settings of baselines and GARNET are available in Appendix D. Our hardware
 300 information is provided in Appendix E.

301 4.1 Robustness of GARNET

302 **Defense on homophilic graphs.** We first evaluate the model robustness on homophilic graphs against
 303 the targeted attack (Nettack) and the non-targeted attack (Metattack). Specifically, Nettack aims to
 304 fool a GNN model to misclassify some target nodes with a few structure (edge) perturbations. The
 305 goal of Metattack is to drop the overall accuracy of the whole test set with a given perturbation ratio
 306 budget (i.e., the number of adversarial edges over the number of total edges). Due to the space limit,
 307 we only show defense results under Nettack and Metattack with 5 perturbed edges per target node
 308 and 20% perturbation ratio, respectively. Results with other perturbation budgets are in Appendix I.

309 Table 1 reports the average accuracy over 10 runs on Cora and Pubmed. It shows that GARNET,
 310 with either a backbone GNN model (GCN or GPRGNN), outperforms defense baselines in terms of
 311 both clean and adversarial accuracy in most cases. We attribute the large accuracy improvement to
 312 GARNET’s strengths in recovering key structures of the clean graph while ignoring the high-rank
 313 adversarial components during graph purification. Moreover, as both TSVD and ProGNN involve
 314 dense matrices during GNN training, they run out of GPU memory even on Pubmed, a graph with
 315 only 20k nodes. In contrast, GARNET is not only robust to adversarial attacks, but also scalable to
 316 large graphs, as empirically shown in Section 4.2.

317 **Defense on heterophilic graphs.** We report the averaged accuracy over 10 runs on heterophilic
 318 graphs in Table 2, which shows that all defense baselines fail to defend GPRGNN on heterophilic
 319 graphs and even degrade the accuracy of the vanilla GPRGNN by a large margin. The reason why
 320 ProGNN performs poorly is that it follows the graph homophily assumption for improving GNN
 321 robustness, which contradicts the property of heterophilic graphs. For the TSVD-based defense
 322 method, the low-rank graph generated by TSVD contains negative edge weights, which degrade the
 323 performance of GPRGNN for adapting its graph filter on heterophilic graphs [31]. Although [29]
 324 have shown GNNGuard can improve model robustness on synthetic heterophilic graphs, our results
 325 indicate that it fails to defend GNN models on realistic heterophilic graphs. We attribute it to that the
 326 quality of graphlet degree vectors used in GNNGuard is degraded by structural perturbations induced
 327 via adversarial attacks. In contrast, GARNET largely recovers the clean graph structure based on
 328 Theorem 3.3 without the assumption on whether adjacent nodes have similar attributes. In other
 329 words, GARNET will produce a heterophilic graph if the underlying clean graph is heterophilic,
 330

Table 1: Averaged node classification accuracy (%) \pm std under targeted attack (Nettack) and non-targeted attack (Metattack) on homophilic graphs — We bold and underline the first and second highest accuracy of each backbone GNN model, respectively. *OOM* means out of memory.

Model	Cora (Nettack)		Cora (Metattack)		Pubmed (Nettack)		Pubmed (Metattack)	
	Clean	Adversarial	Clean	Adversarial	Clean	Adversarial	Clean	Adversarial
GCN-Vanilla	<u>80.96</u> \pm 0.95	55.66 \pm 1.95	81.35 \pm 0.66	56.28 \pm 1.19	87.26 \pm 0.51	66.67 \pm 1.34	87.16 \pm 0.09	77.20 \pm 0.27
GCN-TSVD	72.65 \pm 2.29	60.30 \pm 2.25	73.86 \pm 0.53	62.44 \pm 1.16	87.03 \pm 0.48	<u>79.56</u> \pm 0.48	84.53 \pm 0.08	<u>84.30</u> \pm 0.08
GCN-ProGNN	80.54 \pm 1.21	<u>65.38</u> \pm 1.65	78.56 \pm 0.36	<u>72.28</u> \pm 1.67	88.14 \pm 1.44	71.89 \pm 1.56	84.62 \pm 0.11	83.89 \pm 0.32
GCN-LFR	80.07 \pm 0.95	53.73 \pm 2.17	77.23 \pm 2.61	65.38 \pm 3.71	87.20 \pm 1.24	68.49 \pm 2.44	81.91 \pm 0.26	78.32 \pm 0.69
GCN-GARNET	81.08 \pm 2.05	67.04 \pm 2.05	<u>79.64</u> \pm 0.75	73.89 \pm 0.91	<u>87.96</u> \pm 0.58	86.12 \pm 0.86	<u>85.37</u> \pm 0.20	85.14 \pm 0.23
GPR-Vanilla	83.04 \pm 2.05	62.89 \pm 1.95	83.05 \pm 0.42	74.27 \pm 2.11	<u>90.05</u> \pm 0.73	<u>76.99</u> \pm 1.16	87.35 \pm 0.13	<u>84.18</u> \pm 0.15
GPR-TSVD	81.68 \pm 1.78	63.52 \pm 3.27	81.61 \pm 0.54	<u>78.50</u> \pm 1.20	<i>OOM</i>	<i>OOM</i>	<i>OOM</i>	<i>OOM</i>
GPR-ProGNN	82.04 \pm 1.33	<u>63.74</u> \pm 2.57	82.04 \pm 0.90	76.29 \pm 1.46	<i>OOM</i>	<i>OOM</i>	<i>OOM</i>	<i>OOM</i>
GPR-GARNET	<u>82.77</u> \pm 1.89	71.45 \pm 2.73	<u>82.67</u> \pm 1.89	81.34 \pm 0.79	90.99 \pm 0.52	89.52 \pm 0.45	<u>86.86</u> \pm 0.57	85.69 \pm 0.26

Table 2: Averaged node classification accuracy (%) \pm std on heterophilic graphs — We bold and underline the first and second highest accuracy, respectively. The backbone GNN model is GPRGNN.

Model	Chameleon (Nettack)		Chameleon (Metattack)		Squirrel (Nettack)		Squirrel (Metattack)	
	Clean	Adversarial	Clean	Adversarial	Clean	Adversarial	Clean	Adversarial
Vanilla	<u>71.46</u> \pm 1.92	<u>66.26</u> \pm 1.71	61.36 \pm 1.00	<u>53.20</u> \pm 0.88	<u>41.36</u> \pm 2.87	<u>39.45</u> \pm 2.36	<u>39.51</u> \pm 1.64	<u>35.22</u> \pm 1.20
TSVD	62.12 \pm 3.04	60.37 \pm 2.86	47.29 \pm 1.63	45.12 \pm 1.34	32.98 \pm 2.36	31.20 \pm 1.84	31.36 \pm 1.87	23.91 \pm 1.40
ProGNN	58.80 \pm 1.72	57.07 \pm 1.82	48.39 \pm 0.68	46.69 \pm 0.61	31.81 \pm 1.72	27.27 \pm 1.87	31.64 \pm 2.87	29.36 \pm 3.61
GNNGuard	64.87 \pm 2.62	62.21 \pm 1.94	58.01 \pm 1.57	49.89 \pm 1.34	34.17 \pm 2.33	33.41 \pm 1.82	37.46 \pm 0.56	32.69 \pm 0.59
GARNET	72.89 \pm 2.65	71.83 \pm 2.11	<u>61.11</u> \pm 2.46	59.96 \pm 0.84	44.91 \pm 1.53	43.64 \pm 1.53	43.43 \pm 1.14	41.97 \pm 1.02

331 which is further confirmed in Appendix O. Consequently, GARNET improves accuracy over defense
 332 baselines by up to 10.23% (i.e., 43.64% – 33.41% on Squirrel under Nettack) on heterophilic graphs.

333 **Defense against adaptive attacks.** As GARNET is non-differentiable during kNN graph construction,
 334 it is difficult to optimize a specific loss function for adaptive attack. Instead, we adopt an attack
 335 called LowBlow from [7], which deliberately perturbs low-rank singular components in the graph
 336 spectrum, yet violates the unnoticeable condition (i.e., preserving node degree distribution after
 337 attacking). Since LowBlow has cubic complexity for computing the full set of adjacency eigenpairs,
 338 we only show results on the small graph Cora in Table 3, which indicates GARNET still achieves
 339 the highest adversarial accuracy under LowBlow, while all low-rank defense baselines perform even
 340 worse than vanilla GPRGNN model. The reason lies in that the kNN graph (with a relatively large k)
 341 in GARNET is less vulnerable to the perturbations of weighted spectral embeddings (i.e., low-rank
 342 components of the clean graph) [33], compared to prior low-rank defense methods.

343 **4.2 Scalability of GARNET**

344 To demonstrate the scalability of GARNET, we first compare the run time of GARNET with prior
 345 low-rank defense methods with GPRGNN as the backbone GNN model. As shown in Figure 3, the
 346 TSVD defense method is slower than GARNET since it produces a dense adjacency matrix that
 347 slows down the GNN training. Moreover, ProGNN is extremely slow as it jointly learns the low-rank
 348 graph structure and the robust GNN model, which requires performing TSVD for every epoch. In
 349 contrast, GARNET can efficiently produce a sparse graph for downstream GNN training, leading to
 350 end-to-end runtime speedup over prior methods by up to 14.7 \times . In addition, we further evaluate the
 351 robustness of GARNET on two large datasets: ogbn-arxiv and ogbn-products, under powerful and
 352 scalable attacks proposed by [34]. As we run out of GPU memory when performing the PR-BCD
 353 attack, we choose the more scalable version GR-BCD that has less memory usage. We use GCN as
 354 the backbone model since it outperforms GPRGNN on large graphs. As TSVD and ProGNN run

Table 3: Averaged accuracy (%) \pm std on Cora under Metattack and LowBlow with 20% perturbation ratio. We use GPRGNN as the backbone GNN model.

Model	Metattack	LowBlow
Vanilla	74.27 \pm 2.11	74.77 \pm 0.71
TSVD	78.50 \pm 1.20	26.03 \pm 2.76
ProGNN	76.29 \pm 1.46	69.88 \pm 1.61
GARNET	81.34 \pm 0.79	77.71 \pm 0.95

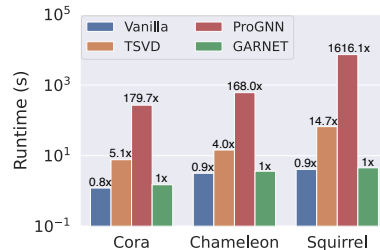
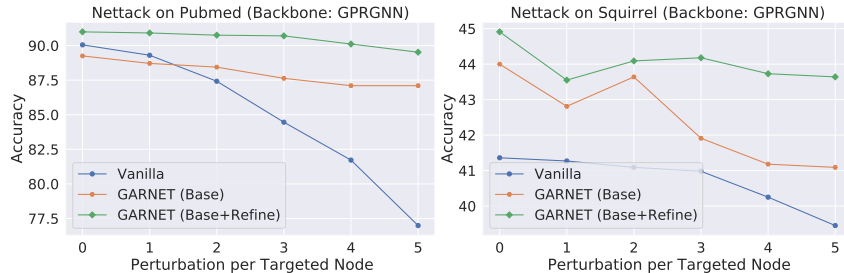


Figure 3: End-to-end runtime comparison of GARNET and baseline methods.

Table 4: Averaged accuracy (%) \pm std under GR-BCD attack.

Model	ogbn-arxiv			ogbn-products		
	Clean	25% Ptb.	50% Ptb.	Clean	25% Ptb.	50% Ptb.
GCN	70.74 \pm 0.26	45.18 \pm 0.25	39.12 \pm 0.27	75.68 \pm 0.20	64.70 \pm 0.43	62.71 \pm 0.44
GNNGuard	68.78 \pm 0.32	<u>47.46</u> \pm 0.11	<u>41.18</u> \pm 0.12	74.82 \pm 0.11	<u>66.76</u> \pm 0.23	<u>63.22</u> \pm 0.26
GCNJaccard	67.77 \pm 0.18	46.27 \pm 0.11	40.84 \pm 0.19	72.95 \pm 0.08	60.90 \pm 0.18	58.84 \pm 0.20
Soft Median GDC	69.75 \pm 0.03	45.31 \pm 0.06	40.11 \pm 0.06	66.31 \pm 0.03	60.59 \pm 0.05	59.73 \pm 0.05
GARNET	<u>69.91</u> \pm 0.29	61.32 \pm 0.20	60.88 \pm 0.13	76.05 \pm 0.19	75.03 \pm 0.14	74.97 \pm 0.24

**Figure 4:** Ablation study of GARNET on graph refinement.

355 out of memory on these two datasets, we choose GNNGuard, GCNJaccard [23], and Soft Median
 356 GDC [34] as baselines. Table 4 shows GARNET achieves comparable clean accuracy compared to
 357 GCN, and drastically improves the adversarial accuracy over defense baselines by up to 16.13%.

358 4.3 Ablation Analysis of GARNET

359 Figure 4 shows the comparison of GARNET results with and without graph refinement. When only
 360 constructing the base graph, GARNET achieves better adversarial accuracy than the vanilla GNN
 361 model, which confirms our Theorem 3.3 that the base graph construction can successfully recover
 362 clean graph edges. The graph refinement step further improves GARNET accuracy ($\sim 2\%$ increase)
 363 since some noncritical or even harmful edges are removed based on PGM. Due to the space limitation,
 364 the ablation studies of GARNET on the kNN graph and edge pruning are available in Appendix G.

365 4.4 Visualization

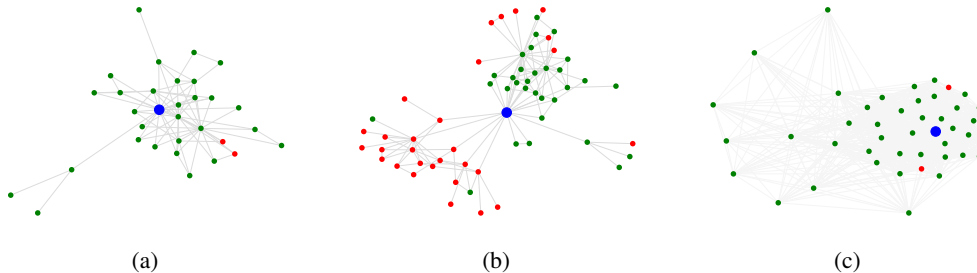


Figure 5: Visualizations on the same target node (marked in blue) as well as its 1-hop and 2-hop neighbors. Neighbor nodes are marked in green if they have the same label as the target node, and red otherwise. (a) clean graph. (b) adversarial graph. (c) adversarial graph purified by GARNET.

366 We visualize the local structure (within 2-hop neighbors) of a target node (randomly picked) on Cora
 367 in Figure 5. By comparing Figures 5(b) and 5(c), it is clear that GARNET effectively removes most
 368 of the adversarial edges induced by Nettack that connect nodes with different labels [8]. As a result, it
 369 is trivial for the backbone GNN model to correctly predict the target node since the surrounding nodes
 370 share the same label as the target node in GARNET graph. This explains why GARNET substantially
 371 improves the adversarial accuracy of GNN models. More visualizations are available in Appendix N.

372 5 Conclusions

373 This work introduces GARNET, a spectral approach to robust and scalable graph neural networks by
 374 combining spectral embedding and the probabilistic graphical model. GARNET first uses weighted
 375 spectral embedding to construct a base graph, which is then refined by pruning uncritical edges based
 376 on the graphical model. Results show that GARNET not only outperforms state-of-the-art defense
 377 models, but also scales to large graphs with millions of nodes. An interesting direction for future
 378 work is to incorporate the node feature information to further boost model robustness.

References

- 379
- 380 [1] William L Hamilton. Graph representation learning. *Synthesis Lectures on Artificial Intelligence*
381 *and Machine Learning*, 14(3):1–159, 2020. 1
- 382 [2] Rex Ying, Ruining He, Kaifeng Chen, Pong Eksombatchai, William L Hamilton, and Jure
383 Leskovec. Graph convolutional neural networks for web-scale recommender systems. In
384 *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery &*
385 *Data Mining*, pages 974–983, 2018. 1
- 386 [3] Sergio Casas, Cole Gulino, Renjie Liao, and Raquel Urtasun. Spaggn: Spatially-aware graph
387 neural networks for relational behavior forecasting from sensor data. In *2020 IEEE International*
388 *Conference on Robotics and Automation (ICRA)*, pages 9491–9497. IEEE, 2020. 1
- 389 [4] Azalia Mirhoseini, Anna Goldie, Mustafa Yazgan, Joe Wenjie Jiang, Ebrahim Songhori, Shen
390 Wang, Young-Joon Lee, Eric Johnson, Omkar Pathak, Azade Nazi, et al. A graph placement
391 methodology for fast chip design. *Nature*, 594(7862):207–212, 2021. 1
- 392 [5] Daniel Zügner, Amir Akbarnejad, and Stephan Günnemann. Adversarial attacks on neural
393 networks for graph data. In *Proceedings of the 24th ACM SIGKDD International Conference*
394 *on Knowledge Discovery & Data Mining*, pages 2847–2856, 2018. 1, 2, 3, 7, 19
- 395 [6] Daniel Zügner and Stephan Günnemann. Adversarial attacks on graph neural networks via meta
396 learning. *arXiv preprint arXiv:1902.08412*, 2019. 1, 2, 3, 7
- 397 [7] Negin Entezari, Saba A Al-Sayouri, Amirali Darvishzadeh, and Evangelos E Papalexakis. All
398 you need is low (rank) defending against adversarial attacks on graphs. In *Proceedings of the*
399 *13th International Conference on Web Search and Data Mining*, pages 169–177, 2020. 1, 3, 4,
400 5, 6, 7, 8, 19
- 401 [8] Wei Jin, Yao Ma, Xiaorui Liu, Xianfeng Tang, Suhang Wang, and Jiliang Tang. Graph
402 structure learning for robust graph neural networks. In *Proceedings of the 26th ACM SIGKDD*
403 *International Conference on Knowledge Discovery & Data Mining*, pages 66–74, 2020. 1, 3, 4,
404 7, 9, 14
- 405 [9] Xiaowen Dong, Dorina Thanou, Michael Rabbat, and Pascal Frossard. Learning graphs from
406 data: A signal representation perspective. *IEEE Signal Processing Magazine*, 36(3):44–63,
407 2019. 2, 3
- 408 [10] Jiong Zhu, Yujun Yan, Lingxiao Zhao, Mark Heimann, Leman Akoglu, and Danai Koutra.
409 Beyond homophily in graph neural networks: Current limitations and effective designs. *arXiv*
410 *preprint arXiv:2006.11468*, 2020. 2, 7, 14, 17, 21
- 411 [11] Weihua Hu, Matthias Fey, Marinka Zitnik, Yuxiao Dong, Hongyu Ren, Bowen Liu, Michele
412 Catasta, and Jure Leskovec. Open graph benchmark: Datasets for machine learning on graphs.
413 *arXiv preprint arXiv:2005.00687*, 2020. 2, 7, 14
- 414 [12] Onureena Banerjee, Laurent El Ghaoui, and Alexandre d’Aspremont. Model selection through
415 sparse maximum likelihood estimation for multivariate gaussian or binary data. *The Journal of*
416 *Machine Learning Research*, 9:485–516, 2008. 3
- 417 [13] Jerome Friedman, Trevor Hastie, and Robert Tibshirani. Sparse inverse covariance estimation
418 with the graphical lasso. *Biostatistics*, 9(3):432–441, 2008. 3, 5
- 419 [14] Xiaowen Dong, Dorina Thanou, Pascal Frossard, and Pierre Vandergheynst. Learning laplacian
420 matrix in smooth graph signal representations. *IEEE Transactions on Signal Processing*, 64
421 (23):6160–6173, 2016. 3
- 422 [15] Hilmi E Egilmez, Eduardo Pavez, and Antonio Ortega. Graph learning from data under laplacian
423 and structural constraints. *IEEE Journal of Selected Topics in Signal Processing*, 11(6):825–841,
424 2017.
- 425 [16] Vassilis Kalofolias and Nathanaël Perraudin. Large scale graph learning from smooth signals.
426 *International Conference on Learning Representations (ICLR 2019)*, 2019.
- 427 [17] Zhuo Feng. Sgl: Spectral graph learning from measurements. *arXiv preprint arXiv:2104.07867*,
428 2021. 3, 6, 22
- 429 [18] Martin Slawski and Matthias Hein. Estimation of positive definite m-matrices and structure
430 learning for attractive gaussian markov random fields. *Linear Algebra and its Applications*, 473:
431 145–179, 2015. 3

- 432 [19] Lichao Sun, Yingtong Dou, Carl Yang, Ji Wang, Philip S Yu, Lifang He, and Bo Li. Adversarial
433 attack and defense on graph data: A survey. *arXiv preprint arXiv:1812.10528*, 2018. 3
- 434 [20] Hanjun Dai, Hui Li, Tian Tian, Xin Huang, Lin Wang, Jun Zhu, and Le Song. Adversarial attack
435 on graph structured data. In *International conference on machine learning*, pages 1115–1124.
436 PMLR, 2018. 3
- 437 [21] Jiong Zhu, Junchen Jin, Michael T Schaub, and Danai Koutra. Improving robustness of graph
438 neural networks with heterophily-inspired designs. *arXiv preprint arXiv:2106.07767*, 2021. 3
- 439 [22] Heng Chang, Yu Rong, Tingyang Xu, Yatao Bian, Shiji Zhou, Xin Wang, Junzhou Huang, and
440 Wenwu Zhu. Not all low-pass filters are robust in graph convolutional networks. In A. Beygelz-
441 imer, Y. Dauphin, P. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information*
442 *Processing Systems*, 2021. URL <https://openreview.net/forum?id=bDdfxLQITtu>. 3, 7
- 443 [23] Huijun Wu, Chen Wang, Yuriy Tyshetskiy, Andrew Docherty, Kai Lu, and Liming Zhu. Ad-
444 versarial examples on graph data: Deep insights into attack and defense. *arXiv preprint*
445 *arXiv:1903.01610*, 2019. 3, 9
- 446 [24] Enyan Dai, Wei Jin, Hui Liu, and Suhang Wang. Towards robust graph neural networks for
447 noisy graphs with sparse labels. In *Proceedings of the Fifteenth ACM International Conference*
448 *on Web Search and Data Mining*, pages 181–191, 2022. 3
- 449 [25] James Baglama and Lothar Reichel. Augmented implicitly restarted lanczos bidiagonalization
450 methods. *SIAM Journal on Scientific Computing*, 27(1):19–42, 2005. 4, 7
- 451 [26] Cho-Jui Hsieh, Matyas A Sustik, Inderjit S Dhillon, and Pradeep Ravikumar. Sparse inverse
452 covariance matrix estimation using quadratic approximation. *arXiv preprint arXiv:1306.3212*,
453 2013. 5
- 454 [27] Cho-Jui Hsieh, Mátyás A Sustik, Inderjit S Dhillon, Pradeep Ravikumar, et al. Quic: quadratic
455 approximation for sparse inverse covariance estimation. *J. Mach. Learn. Res.*, 15(1):2911–2947,
456 2014. 5
- 457 [28] Yu A Malkov and Dmitry A Yashunin. Efficient and robust approximate nearest neighbor search
458 using hierarchical navigable small world graphs. *IEEE transactions on pattern analysis and*
459 *machine intelligence*, 42(4):824–836, 2018. 6, 7, 15
- 460 [29] Xiang Zhang and Marinka Zitnik. Gnn-guard: Defending graph neural networks against
461 adversarial attacks. *arXiv preprint arXiv:2006.08149*, 2020. 7
- 462 [30] Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional
463 networks. *arXiv preprint arXiv:1609.02907*, 2016. 7, 17
- 464 [31] Eli Chien, Jianhao Peng, Pan Li, and Olgica Milenkovic. Adaptive universal generalized
465 pagerank graph neural network. In *International Conference on Learning Representations*,
466 2021. URL <https://openreview.net/forum?id=n6jl7fLxrP>. 7, 14, 17
- 467 [32] Derek Lim, Felix Matthew Hohne, Xiuyu Li, Sijia Linda Huang, Vaishnavi Gupta, Omkar Prasad
468 Bhalerao, and Ser-Nam Lim. Large scale learning on non-homophilous graphs: New bench-
469 marks and strong simple methods. In *Advances in Neural Information Processing Systems*,
470 2021. 7
- 471 [33] Yizhen Wang, Somesh Jha, and Kamalika Chaudhuri. Analyzing the robustness of nearest
472 neighbors to adversarial examples. In *International Conference on Machine Learning*, pages
473 5133–5142. PMLR, 2018. 8
- 474 [34] Simon Geisler, Tobias Schmidt, Hakan Şirin, Daniel Zügner, Aleksandar Bojchevski, and
475 Stephan Günnemann. Robustness of graph neural networks at scale. *Advances in Neural*
476 *Information Processing Systems*, 34, 2021. 8, 9, 14
- 477 [35] Gene H Golub and Charles F Van Loan. *Matrix computations*, volume 3. JHU press, 2013. 13
- 478 [36] Zhilin Yang, William Cohen, and Ruslan Salakhudinov. Revisiting semi-supervised learning
479 with graph embeddings. In *International conference on machine learning*, pages 40–48. PMLR,
480 2016. 14
- 481 [37] Benedek Rozemberczki, Carl Allen, and Rik Sarkar. Multi-scale attributed node embedding.
482 *Journal of Complex Networks*, 9(2):cnab014, 2021. 14
- 483 [38] Yaxin Li, Wei Jin, Han Xu, and Jiliang Tang. Deeprobust: A pytorch library for adversarial
484 attacks and defenses. *arXiv preprint arXiv:2005.06149*, 2020. 14

- 485 [39] Ulrike Von Luxburg. A tutorial on spectral clustering. *Statistics and computing*, 17(4):395–416,
486 2007. 15
- 487 [40] Chenhui Deng, Zhiqiang Zhao, Yongyu Wang, Zhiru Zhang, and Zhuo Feng. Graphzoom:
488 A multi-level spectral approach for accurate and scalable graph embedding. In *International*
489 *Conference on Learning Representations*, 2020. URL [https://openreview.net/forum?](https://openreview.net/forum?id=r11G00EKDH)
490 [id=r11G00EKDH](https://openreview.net/forum?id=r11G00EKDH). 19
- 491 [41] Wuxinlin Cheng, Chenhui Deng, Zhiqiang Zhao, Yaohui Cai, Zhiru Zhang, and Zhuo Feng.
492 Spade: A spectral method for black-box adversarial robustness evaluation. In *International*
493 *Conference on Machine Learning*, pages 1814–1824. PMLR, 2021. 19

494 **A Proof for Proposition 3.2**

495 *Proof.* As the graph is undirected, we can perform eigendecomposition on both A_{norm} and L_{norm}
 496 to obtain their real eigenvalues and the corresponding eigenvectors. Let λ_i , $\hat{\lambda}_i$, and σ_i , $i = 1, 2, \dots, r$
 497 denote the r smallest eigenvalues of L_{norm} , r largest eigenvalues of A_{norm} , and r largest singular
 498 values of A_{norm} , respectively. Since $A_{norm} = I - L_{norm}$, A_{norm} and L_{norm} share the same set of
 499 eigenvectors while their eigenvalues satisfy: $\hat{\lambda}_i = 1 - \lambda_i$, $i = 1, 2, \dots, r$. Moreover, since we assume
 500 that the r largest magnitude eigenvalues of A_{norm} are non-negative, we have $\sigma_i = |\hat{\lambda}_i| = \hat{\lambda}_i$, $i =$
 501 $1, 2, \dots, r$. Thus, we have:

$$\begin{aligned} VV^T &= [v_1, \dots, v_r] \begin{bmatrix} |1 - \lambda_1| & & \\ & \ddots & \\ & & |1 - \lambda_r| \end{bmatrix} [v_1, \dots, v_r]^T \\ &= [v_1, \dots, v_r] \begin{bmatrix} |\hat{\lambda}_1| & & \\ & \ddots & \\ & & |\hat{\lambda}_r| \end{bmatrix} [v_1, \dots, v_r]^T \\ &= [v_1, \dots, v_r] \begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_r \end{bmatrix} [v_1, \dots, v_r]^T \\ &= \hat{A} \end{aligned}$$

502

□

 503 **B Proof for Theorem 3.3**

504 *Proof.* Since the weighted embedding matrix V is defined as $V \stackrel{\text{def}}{=} [\sqrt{|1 - \lambda_1|}v_1, \dots, \sqrt{|1 - \lambda_r|}v_r]$,
 505 where $\lambda_1, \lambda_2, \dots, \lambda_r$ and v_1, v_2, \dots, v_r are the top r smallest eigenvalues and the corresponding
 506 eigenvectors of normalized graph Laplacian matrix $L_{norm} = I - D^{-\frac{1}{2}}AD^{-\frac{1}{2}}$, we have:

$$\begin{aligned} \sum_{(i,j) \in \mathcal{E}} \|V_i - V_j\|_2^2 &= \sum_{k=1}^r \sum_{(i,j) \in \mathcal{E}} |1 - \lambda_k| (v_{k,i} - v_{k,j})^2 \\ &= \sum_{k=1}^r |1 - \lambda_k| \sum_{(i,j) \in \mathcal{E}} (v_{k,i} - v_{k,j})^2 \\ &= \sum_{k=1}^r |1 - \lambda_k| v_k^T L_{norm} v_k \\ &= \sum_{k=1}^r (1 - \lambda_k) \lambda_k \\ &\leq \sum_{k=1}^r 0.25 \\ &= 0.25r \end{aligned}$$

507

□

508 The fourth equation above is based on Courant-Fischer Theorem [35] with the assumption that $\lambda_r \leq 1$
 509 and the Laplacian eigenvectors are normalized (i.e., $\|v_k\|_2 = 1, \forall k = 1, \dots, r$). The inequality is
 510 derived by arithmetic mean-geometric mean (AM-GM) inequality.

Table 5: Statistics of datasets used in our experiments.

Dataset	Type	Homophily Score	Nodes	Edges	Classes	Features
Cora	Homophily	0.80	2,485	5,069	7	1,433
Pubmed	Homophily	0.80	19,717	44,324	3	500
Chameleon	Heterophily	0.23	2,277	62,792	5	2,325
Squirrel	Heterophily	0.22	5,201	396,846	5	2,089
ogbn-arxiv	Homophily	0.66	169,343	1,166,243	40	128
ogbn-products	Homophily	0.81	2,449,029	61,859,140	47	100

511 C Dataset Details

512 Table 5 shows the statistics of the datasets used in our experiments. We follow Zhu et al. [10] to
 513 compute the homophily score per dataset (lower score means more heterophilic). As in Jin et al. [8],
 514 we extract the largest connected components of the original Cora and Pubmed datasets [36] for the
 515 adversarial evaluation, with the same train/validation/test split. For Chameleon and Squirrel [37], we
 516 keep the same split setting as Chien et al. [31]. Finally, we follow the split setting of Open Graph
 517 Benchmark (OGB) [11] on ogbn-arxiv and ogbn-products. Note that all data used in our experiments
 518 do not contain personally identifiable information or offensive content.

519 In addition, we follow Jin et al. [8] for the selection of target nodes on Cora and Pubmed under
 520 Nettack. For the Chameleon and Squirrel datasets under Nettack, we choose target nodes that have
 521 degrees within the range of [20, 50] and [20, 140], respectively. In regard to non-targeted attacks (i.e.,
 522 Metattack), we choose nodes in the test set as target nodes for all datasets. We implement all the
 523 adversarial attacks based on the DeepRobust library [38].

524 D Hyperparameters Settings

525 D.1 Backbone GNN Models

526 **GCN.** We choose the GCN hyperparameters based on the DeepRobust library [38].

527 **GPRGNN.** We follow the hyperparameter settings provided at github.com/jianhao2016/GPRGNN
 528 with slightly different dropout rates (chosen from 0.3, 0.5, 0.7) and learning rates (chosen from
 529 0.01, 0.05, 0.1). Specifically, we provide the complete choices of dropout rates and learning rates
 530 across all datasets and attack settings below:

- 531 • Cora-Nettack: dropout of 0.5 and learning rate of 0.01.
- 532 • Cora-Metattack: dropout of 0.5 and learning rate of 0.01.
- 533 • Pubmed-Nettack: dropout of 0.5 and learning rate of 0.01.
- 534 • Pubmed-Metattack: dropout of 0.5 and learning rate of 0.01.
- 535 • Chameleon-Nettack: dropout of 0.5 and learning rate of 0.05.
- 536 • Chameleon-Metattack: dropout of 0.3 and learning rate of 0.05.
- 537 • Squirrel-Nettack: dropout of 0.5 and learning rate of 0.1.
- 538 • Squirrel-Metattack: dropout of 0.5 and learning rate of 0.1.

539 **H2GCN.** We train a three-layer model in full batch, with a learning rate of 0.01, dropout of 0.5,
 540 hidden dimension of 64, and 300 epochs for both Chameleon and Squirrel datasets.

541 D.2 Defense Baselines

542 **TSVD.** We use the same r eigenvectors in TSVD as those used in GARNET, which is shown in
 543 Table 6.

544 **GCNJaccard.** We choose the GCNJaccard hyperparameters based on the DeepRobust library [38].

545 **GNNGuard.** We set edge pruning threshold (the only hyperparameter in GNNGuard) to be $P_0 = 0.1$.

546 **Soft Median GDC.** We strictly follow the hyperparameter setting suggested by Geisler et al. [34]. In
 547 particular, we choose temperature T of 5.0 for soft median, α of 0.1 (0.15) and k of 64 (32) for GDC
 548 on ogbn-arxiv (ogbn-products).

549 **ProGNN.** We find out its performance is very sensitive to hyperparameters. Thus we strictly follow the
 550 tuned hyperparameters available at github.com/ChandlerBang/Pro-GNN/scripts. As GCN-ProGNN
 551 training is very slow on Pubmed (estimated time is 30 days for 10 runs), we follow the suggestion
 552 from ProGNN authors to replace “svd” with “truncated svd” in the ProGNN implementation.

553 D.3 GARNET

Table 6: Summary of hyperparameters in GARNET— We denote the number of eigenpairs for spectral embedding, the number of nearest neighbors for base graph construction, and the threshold for edge pruning by r , k , and γ , respectively.

Dataset	r	k	γ
Cora-Nettack	50	30	0.003
Cora-Metattack	50	30	0.003
Pubmed-Nettack	50	50	0.003
Pubmed-Metattack	50	50	0.003
Chameleon-Nettack	50	50	0.003
Chameleon-Metattack	50	50	0.003
Squirrel-Nettack	50	50	0.003
Squirrel-Metattack	50	50	0.003
ogbn-arxiv-GRBCD	500	50	0.003
ogbn-products-GRBCD	500	50	0.003

554 We show the hyperparameters of GARNET on different datasets under Nettack (1 perturbation per
 555 node), Metattack (10% perturbation ratio), and GR-BCD (25% perturbation ratio) in Table 6. Note
 556 that we provide our strategy of choosing r in Appendix F, which avoids conducting hyperparameter
 557 tuning on r per dataset. Besides, we set the prior data variance σ^2 to be 1000 for all graphs. In
 558 addition, we run all GNN training with a full batch way.

559 E Hardware Information

560 We conduct all experiments on a Linux machine with an Intel Xeon Gold 5218 CPU (8 cores @
 561 2.30GHz) CPU, 8 NVIDIA RTX 2080 Ti GPU (11 GB memory per GPU), and 1 RTX A6000 GPU
 562 (48 GB memory).

563 F Complexity Analysis of GARNET

564 F.1 Time Complexity – Choice of r

565 We choose r based on the number of classes per dataset, which depends on the downstream task
 566 rather than number of nodes in the graph. Specifically, suppose λ_r is the r -th largest eigenvalue, an
 567 appropriate r is chosen if there is a large gap between λ_r and λ_{r+1} (i.e., a large eigengap) in the
 568 graph spectrum. According to [39], the eigengap is highly related to the number of clusters in the
 569 graph. In this work, we approximate r by $r \approx 10c$ to cover the large eigengap, where c denotes
 570 the number of classes/clusters. As shown in Tables 5 and 6, the number of classes in small (large)
 571 graphs is around 5 (50), so we use $r = 50$ ($r = 500$) in experiments. As a result, GARNET has the
 572 near-linear time complexity $O(r(|E| + k|V|) + |V|\log|V|) = O(c(|E| + k|V|) + |V|\log|V|)$.

573 F.2 Space Complexity

574 GARNET involves forming a sparse kNN graph by building hierarchical navigable small world
 575 (HNSW) graphs [28] that contain $O(|V|\log|V|)$ nodes in total and each node connects to a fixed
 576 number of neighbors. Thus, the space complexity of storing the HNSW graphs is $O(|V|\log|V|)$. In
 577 addition, GARNET also needs to store the input adversarial graph and the produced kNN graph. As
 578 a result, the total space complexity of GARNET is $O(|V|(\log|V| + k) + |\mathcal{E}|)$, where $|V|$ and $|\mathcal{E}|$
 579 denote the number of nodes and edges in the adversarial graph, respectively, and k is the averaged
 580 node degree in the kNN graph.

581 Apart from the complexity analysis, we further provide the algorithms of GARNET and TSVD below
 582 for comparison.

Algorithm 1: GARNET based adversarial defense (this work)

Input: Adversarial graph \mathcal{G}_{adv} ; node feature matrix $X \in R^{n \times d}$; prior data variance σ^2 ; truncated svd rank r ; kNN graph k ; threshold for edge pruning γ ; a GNN model for defense.

Output: Node embedding matrix $Z \in R^{n \times c}$

```

1  $M, S = eig_s(\mathcal{G}_{adv}, r)$ ;
2  $V = M \sqrt{|I - S|}$ ;
3  $\mathcal{G}_{base} = kNN\_graph(V, k)$ ;
583 4  $M', S' = eig_s(\mathcal{G}_{base}, r)$ ;
5  $U = M' / \sqrt{S' + I / \sigma^2}$ ;
6 for  $e_{i,j} \in \mathcal{G}_{base}$  do
7   if  $\frac{\|U_i - U_j\|_2^2}{\|V_i - V_j\|_2^2} < \gamma$  then
8     Prune  $e_{i,j}$  from  $\mathcal{G}_{base}$ ;
9   end
10 end
11  $Z = GNN(\mathcal{G}'_{base}, X)$ ;

```

Algorithm 2: Truncated SVD based adversarial defense (prior work)

Input: Adversarial graph \mathcal{G}_{adv} ; node feature matrix $X \in R^{n \times d}$; truncated svd rank r ; a GNN model for defense.

Output: Node embedding matrix $Z \in R^{n \times c}$

```

584 12  $U, S, V = TSV D(\mathcal{G}_{adv}, r)$ ;
13  $A_{\mathcal{G}_{tsvd}} = USV^T$ ;
14  $Z = GNN(\mathcal{G}_{tsvd}, X)$ ;

```

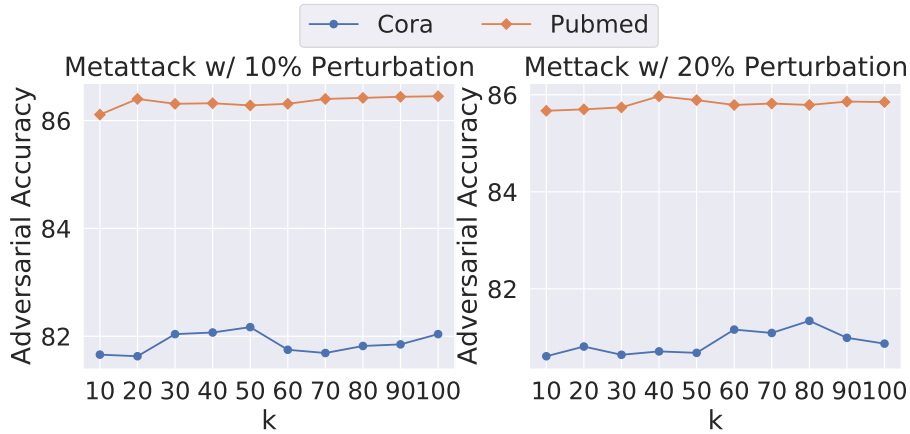
G Ablation Study**G.1 Choice of k for kNN Graph Construction**

Figure 6: Ablation study of GARNET on k for kNN graph construction.

587 To evaluate the sensitivity of GARNET to k nearest-neighbor (kNN) graph construction, we evaluate
588 the adversarial accuracy of GARNET with different k values for constructing kNN graphs. Figure

6 shows that the accuracy of GARNET does not change too much when varying k value within the range of $[10, 100]$, indicating a relatively large k (e.g., $k \geq 10$) can enable the kNN graph to incorporate most of edges in the underlying clean graph. Consequently, the performance of GARNET is relatively robust to the choice of k for kNN graph construction. As the peak performance is typically achieved in $[30, 80]$, we recommend choosing $k = 30 \sim 80$ for building the kNN graph in practice.

G.2 Choice of γ for edge pruning

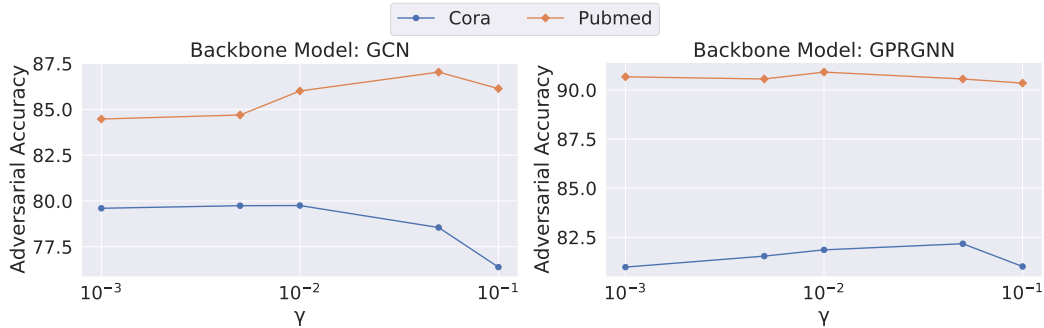


Figure 7: Ablation study of GARNET on threshold γ for edge pruning.

Apart from the choice of k for kNN graph construction, another critical hyperparameter of GARNET is the threshold γ that determines whether an edge should be pruned in the base graph. Thus, we further evaluate the effect of γ on the performance of GARNET. Specifically, we pick γ in the set of $\{0.001, 0.005, 0.01, 0.05, 0.1\}$ and evaluate the corresponding adversarial accuracy of GARNET under Nettack with 1 perturbation per target node. As shown in Figure 7, the improper choice of γ may degrade the adversarial accuracy of GARNET by 3%. However, picking γ around 0.01 can always achieve a reasonable adversarial accuracy on different datasets with different backbone GNN models. As a result, we recommend choosing γ in the range of $[0.005, 0.05]$ in practice. In addition, Figure 7 also shows that GPRGNN-GARNET is more robust to the changes of γ than GCN-GARNET, which is due to the reason that the adaptive graph filter in GPRGNN can adapt to the graph structure after edge pruning.

H Backbone GNN Models for Defense

As GARNET can be integrated with any existing GNN models to improve their adversarial accuracy, we choose two popular GNN models as the backbone model in our experiments: GCN and GPRGNN [30, 31]. As the GCN model implicitly assumes the underlying graph is homophilic, it performs poorly on heterophilic graphs [10]. In contrast, GPRGNN can work on both homophilic and heterophilic datasets, due to its learned graph filter that can adapt to the homophily/heterophily property of the underlying graph. Thus, we choose GPRGNN as the backbone model for evaluation on heterophilic datasets. In addition, we also show the defense results with the H2GCN [10] as backbone model in Appendix J.

I Defense Results with Various Perturbation Budgets

We provide additional defense results under Nettack and Metattack with various perturbations in Tables 7 and 8 respectively. The results indicate that GARNET outperforms prior defense methods in most cases.

J Defense on H2GCN

We provide the results of combining GARNET with H2GCN [10] on heterophilic graphs in Table 9, which shows that GARNET achieves the highest accuracy in most cases and improves the accuracy on all perturbed graphs by a large margin compared to the vanilla H2GCN as well as H2GCN-TSVD.

Table 7: Averaged node classification accuracy (%) \pm std under targeted attack (Nettack) with different perturbation ratio — We denote the evaluated dataset by its name with the number of perturbations (e.g., Cora-0 means the clean Cora graph and Cora-1 denotes there is 1 adversarial edge perturbation per target node). As GCN is not designed for heterophilic graphs, we only show results of defense methods with GPRGNN as the backbone model on Chameleon and Squirrel. We bold and underline the first and second highest accuracy of each backbone GNN model, respectively. *OOM* means out of memory.

Dataset	GCN				GPRGNN			
	Vanilla	TSVD	ProGNN	GARNET	Vanilla	TSVD	ProGNN	GARNET
Cora-0	<u>80.96</u> \pm 0.95	72.65 \pm 2.29	80.54 \pm 1.21	81.08 \pm 2.05	83.04 \pm 2.05	81.68 \pm 1.78	82.04 \pm 1.33	<u>82.77</u> \pm 1.89
Cora-1	70.06 \pm 0.81	71.36 \pm 1.63	81.65 \pm 0.59	79.75 \pm 2.35	<u>81.68</u> \pm 2.18	79.36 \pm 2.23	80.56 \pm 1.71	82.17 \pm 1.95
Cora-2	68.60 \pm 1.81	70.66 \pm 2.76	79.83 \pm 1.10	<u>79.69</u> \pm 1.50	74.34 \pm 2.41	<u>76.26</u> \pm 2.34	76.12 \pm 2.43	78.55 \pm 2.11
Cora-3	65.04 \pm 3.31	68.20 \pm 1.93	<u>72.08</u> \pm 1.20	74.42 \pm 2.06	70.96 \pm 2.00	70.90 \pm 3.89	<u>73.74</u> \pm 2.73	79.40 \pm 1.35
Cora-4	61.69 \pm 1.48	65.34 \pm 3.46	<u>67.83</u> \pm 1.87	69.60 \pm 2.67	65.90 \pm 1.61	65.51 \pm 3.27	<u>68.94</u> \pm 3.25	72.77 \pm 2.16
Cora-5	55.66 \pm 1.95	60.30 \pm 2.25	<u>65.38</u> \pm 1.65	67.04 \pm 2.05	62.89 \pm 1.95	63.52 \pm 3.27	<u>63.74</u> \pm 2.57	71.45 \pm 2.73
Pubmed-0	87.26 \pm 0.51	87.03 \pm 0.48	88.14 \pm 1.44	<u>87.96</u> \pm 0.58	<u>90.05</u> \pm 0.73	<i>OOM</i>	<i>OOM</i>	90.99 \pm 0.52
Pubmed-1	84.29 \pm 0.68	<u>86.46</u> \pm 0.28	85.75 \pm 1.23	87.03 \pm 0.68	<u>89.30</u> \pm 0.54	<i>OOM</i>	<i>OOM</i>	90.91 \pm 0.47
Pubmed-2	82.17 \pm 0.67	<u>83.68</u> \pm 0.46	81.23 \pm 1.21	86.92 \pm 0.45	<u>87.42</u> \pm 0.28	<i>OOM</i>	<i>OOM</i>	90.75 \pm 0.55
Pubmed-3	81.13 \pm 0.53	<u>81.34</u> \pm 0.68	80.65 \pm 1.39	86.50 \pm 0.45	<u>84.46</u> \pm 0.53	<i>OOM</i>	<i>OOM</i>	90.70 \pm 0.37
Pubmed-4	75.48 \pm 0.52	<u>82.41</u> \pm 0.54	78.46 \pm 1.11	86.44 \pm 0.64	<u>81.72</u> \pm 0.72	<i>OOM</i>	<i>OOM</i>	90.11 \pm 0.57
Pubmed-5	66.67 \pm 1.34	<u>79.56</u> \pm 0.48	71.89 \pm 1.56	86.12 \pm 0.86	<u>76.99</u> \pm 1.16	<i>OOM</i>	<i>OOM</i>	89.52 \pm 0.45
Chameleon-0					<u>71.46</u> \pm 1.92	62.12 \pm 3.04	58.80 \pm 1.72	72.89 \pm 2.65
Chameleon-1					<u>71.02</u> \pm 1.57	61.34 \pm 2.93	58.05 \pm 1.90	72.68 \pm 1.89
Chameleon-2					<u>70.71</u> \pm 1.12	61.09 \pm 2.80	57.44 \pm 1.67	72.20 \pm 2.31
Chameleon-3					<u>70.30</u> \pm 1.28	60.98 \pm 2.82	57.19 \pm 1.83	72.17 \pm 2.07
Chameleon-4					<u>69.87</u> \pm 1.29	60.85 \pm 3.31	57.44 \pm 1.63	72.06 \pm 2.94
Chameleon-5					<u>66.26</u> \pm 1.71	60.37 \pm 2.86	57.07 \pm 1.82	71.83 \pm 2.11
Squirrel-0					<u>41.36</u> \pm 2.87	32.98 \pm 2.36	31.81 \pm 1.72	44.91 \pm 1.53
Squirrel-1					<u>41.27</u> \pm 3.16	32.63 \pm 0.87	30.54 \pm 2.45	43.55 \pm 1.79
Squirrel-2					<u>41.09</u> \pm 2.14	32.05 \pm 1.05	30.73 \pm 2.13	44.09 \pm 2.35
Squirrel-3					<u>40.98</u> \pm 2.72	32.00 \pm 1.66	30.25 \pm 1.98	44.18 \pm 2.26
Squirrel-4					<u>40.25</u> \pm 2.82	31.45 \pm 1.38	29.09 \pm 2.33	43.73 \pm 1.62
Squirrel-5					<u>39.45</u> \pm 2.36	31.20 \pm 1.84	27.27 \pm 1.87	43.64 \pm 1.53

Table 8: Averaged node classification accuracy (%) \pm std under non-targeted attack (Metattack) with different perturbation ratio — We denote the evaluated dataset by its name with the perturbation ratio (e.g., Cora-0 means the clean Cora graph and Cora-10 denotes there are 10% adversarial edges). As GCN is not designed for heterophilic graphs, we only show results of defense methods with GPRGNN as the backbone model on Chameleon and Squirrel. We bold and underline the first and second highest accuracy of each backbone GNN model, respectively. *OOM* means out of memory.

Dataset	GCN				GPRGNN			
	Vanilla	TSVD	ProGNN	GARNET	Vanilla	TSVD	ProGNN	GARNET
Cora-0	81.35 \pm 0.66	73.86 \pm 0.53	78.56 \pm 0.36	<u>79.64</u> \pm 0.75	83.05 \pm 0.42	81.61 \pm 0.54	82.04 \pm 0.90	<u>82.67</u> \pm 1.89
Cora-10	69.50 \pm 1.46	69.45 \pm 0.69	77.90 \pm 0.69	<u>77.78</u> \pm 0.53	80.37 \pm 0.65	<u>81.08</u> \pm 0.52	80.31 \pm 1.23	82.17 \pm 0.69
Cora-20	56.28 \pm 1.19	62.44 \pm 1.16	<u>72.28</u> \pm 1.67	73.89 \pm 0.91	74.27 \pm 2.11	<u>78.50</u> \pm 1.20	76.29 \pm 1.46	81.34 \pm 0.79
Pubmed-0	87.16 \pm 0.09	84.53 \pm 0.08	84.62 \pm 0.11	<u>85.37</u> \pm 0.20	87.35 \pm 0.13	<i>OOM</i>	<i>OOM</i>	<u>86.86</u> \pm 0.57
Pubmed-10	81.16 \pm 0.13	<u>84.56</u> \pm 0.10	84.09 \pm 0.12	85.22 \pm 0.13	<u>85.52</u> \pm 0.14	<i>OOM</i>	<i>OOM</i>	86.24 \pm 0.20
Pubmed-20	77.20 \pm 0.27	<u>84.30</u> \pm 0.08	83.89 \pm 0.32	85.14 \pm 0.23	<u>84.18</u> \pm 0.15	<i>OOM</i>	<i>OOM</i>	85.69 \pm 0.26
Chameleon-0					61.36 \pm 1.00	47.29 \pm 1.63	48.39 \pm 0.68	<u>61.11</u> \pm 2.46
Chameleon-10					<u>57.55</u> \pm 1.26	47.07 \pm 1.21	47.80 \pm 0.91	60.96 \pm 1.22
Chameleon-20					<u>53.20</u> \pm 0.88	45.12 \pm 1.34	46.69 \pm 0.61	59.96 \pm 0.84
Squirrel-0					<u>39.51</u> \pm 1.64	31.36 \pm 1.87	31.64 \pm 2.87	43.43 \pm 1.14
Squirrel-10					<u>38.27</u> \pm 0.83	28.25 \pm 1.66	30.33 \pm 3.29	42.62 \pm 1.09
Squirrel-20					<u>35.22</u> \pm 1.20	23.91 \pm 1.40	29.36 \pm 3.61	41.97 \pm 1.02

Table 9: Averaged node classification accuracy (%) \pm std on heterophilic graphs — We bold and underline the first and second highest accuracy, respectively. The backbone GNN model is H2GCN.

Model	Chameleon (Nettack)		Chameleon (Metattack)		Squirrel (Nettack)		Squirrel (Metattack)	
	Clean	Adversarial	Clean	Adversarial	Clean	Adversarial	Clean	Adversarial
Vanilla	<u>78.43</u> \pm 2.09	62.20 \pm 1.99	68.45 \pm 0.57	52.73 \pm 1.72	55.36 \pm 2.91	29.55 \pm 3.09	61.23 \pm 0.71	<u>44.84</u> \pm 0.89
TSVD	67.07 \pm 1.15	<u>63.17</u> \pm 1.61	61.75 \pm 1.09	<u>54.06</u> \pm 1.66	32.45 \pm 1.87	<u>31.64</u> \pm 2.09	46.66 \pm 1.71	<u>40.56</u> \pm 1.41
GARNET	78.78 \pm 1.84	76.10 \pm 1.92	<u>66.63</u> \pm 1.05	61.12 \pm 0.59	<u>54.09</u> \pm 1.73	53.27 \pm 1.50	<u>59.67</u> \pm 0.83	50.08 \pm 1.92

624 The results further confirm that GARNET is able to improve robustness of different backbone GNN
 625 models.

626 K Broader Impact

627 Zügner et al. [5] have shown that graph adversarial attacks can drastically degrade the performance
 628 of GNN models for downstream applications. For instance, an attacker can attack a GNN-based
 629 recommender system on Facebook social network or Amazon co-purchasing network, via creating a
 630 fake account and make some connections to other users or items. Those connections can be viewed
 631 as adversarial edges in the graph. As a result, the attacker can deliberately enforce a GNN model to
 632 recommend some irrelevant or even harmful contents to other users. Thus, improving adversarial
 633 robustness of GNN models has the potential for positive societal benefit.

634 We hope that this paper provides insight on the robustness and scalability limitations of prior defense
 635 methods. Moreover, we believe that the proposed GARNET can largely overcome these two limi-
 636 tations and produce a robust GNN model against adversarial attacks on large-scale graph datasets.
 637 Nevertheless, we have to admit that GARNET may potentially provide the attacker with some hints
 638 about developing a even more powerful and scalable adversarial attack than all existing attacks, which
 639 is a possible negative consequence.

640 L Discussion on Node Features

641 L.1 Graph Construction from Node Features

Table 10: Averaged node classification accuracy (%) \pm std under targeted attack (Nettack) and non-targeted attack (Metattack) on Cora and Pubmed — We bold and underline the first and second highest accuracy, respectively. “NodeFeat” denotes the graph constructed from node features is used for GNN training.

Model	Cora (Nettack)		Cora (Metattack)		Pubmed (Nettack)		Pubmed (Metattack)	
	Clean	Adversarial	Clean	Adversarial	Clean	Adversarial	Clean	Adversarial
GCN-Vanilla	<u>80.96</u> \pm 0.95	<u>55.66</u> \pm 1.95	81.35 \pm 0.66	56.28 \pm 1.19	<u>87.26</u> \pm 0.51	66.67 \pm 1.34	87.16 \pm 0.09	77.20 \pm 0.27
GCN-NodeFeat	52.65 \pm 2.69	52.65 \pm 2.69	56.44 \pm 1.04	<u>56.44</u> \pm 1.04	83.01 \pm 0.99	<u>83.01</u> \pm 0.99	78.66 \pm 0.15	<u>78.66</u> \pm 0.15
GCN-GARNET	81.08 \pm 2.05	67.04 \pm 2.05	<u>79.64</u> \pm 0.75	73.89 \pm 0.91	87.96 \pm 0.58	86.12 \pm 0.86	<u>85.37</u> \pm 0.20	85.14 \pm 0.23

642 As GARNET purifies the adversarial graph by building a kNN graph based on dominant singular
 643 components, a natural question is whether the kNN graph constructed from node features can also
 644 achieve similar performance. We answer this question by comparing the results of GARNET graph
 645 and the node feature graph in Table 10. Note that the clean and adversarial accuracy are the same on
 646 the graph constructed from node features, since node features are unchanged after graph adversarial
 647 attack. Besides, we only show results on homophilic graphs as the kNN graph constructed from node
 648 features naturally falls into this category. Table 10 shows that the node feature graph performs much
 649 worse than GARNET graph. This further confirms that the method proposed in this work is critical to
 650 improve the robustness of GNN models.

651 L.2 Defense Against Node Feature Attack

652 GARNET can be extended to handle node feature attack, although this paper mainly focuses on
 653 defending against graph structure attack, which we believe is more challenging than defending node
 654 feature attack due to the discrete nature. Specifically, we can perform TSVD to obtain the low-rank
 655 approximation of the node feature matrix, which can remove high-rank adversarial components
 656 in node features [7]. The low-rank feature matrix is then concatenated to the weighted spectral
 657 embeddings to produce the kNN base graph. In this way, the downstream GNN model will be able to
 658 aggregate neighbors whose features are less perturbed during message passing.

659 M Run Time on OGB Graphs

660 Apart from the runtime comparison of GARNET and defense baselines on small graphs in Figure 3,
 661 we further evaluate the run time of GARNET on large (OGB) graphs. Concretely, the end-to-end run
 662 time of GARNET is 18 mins and 2 hours on ogbn-arxiv and ogbn-products, respectively, which is
 663 $3\times$ faster than the most competitive baseline GNNGuard that takes more than 1 hour on ogbn-arxiv
 664 and 8 hours on ogbn-products. One way to further accelerate GARNET is to leverage prior work on
 665 accelerating spectral embedding [40, 41]. We leave this to our future work.

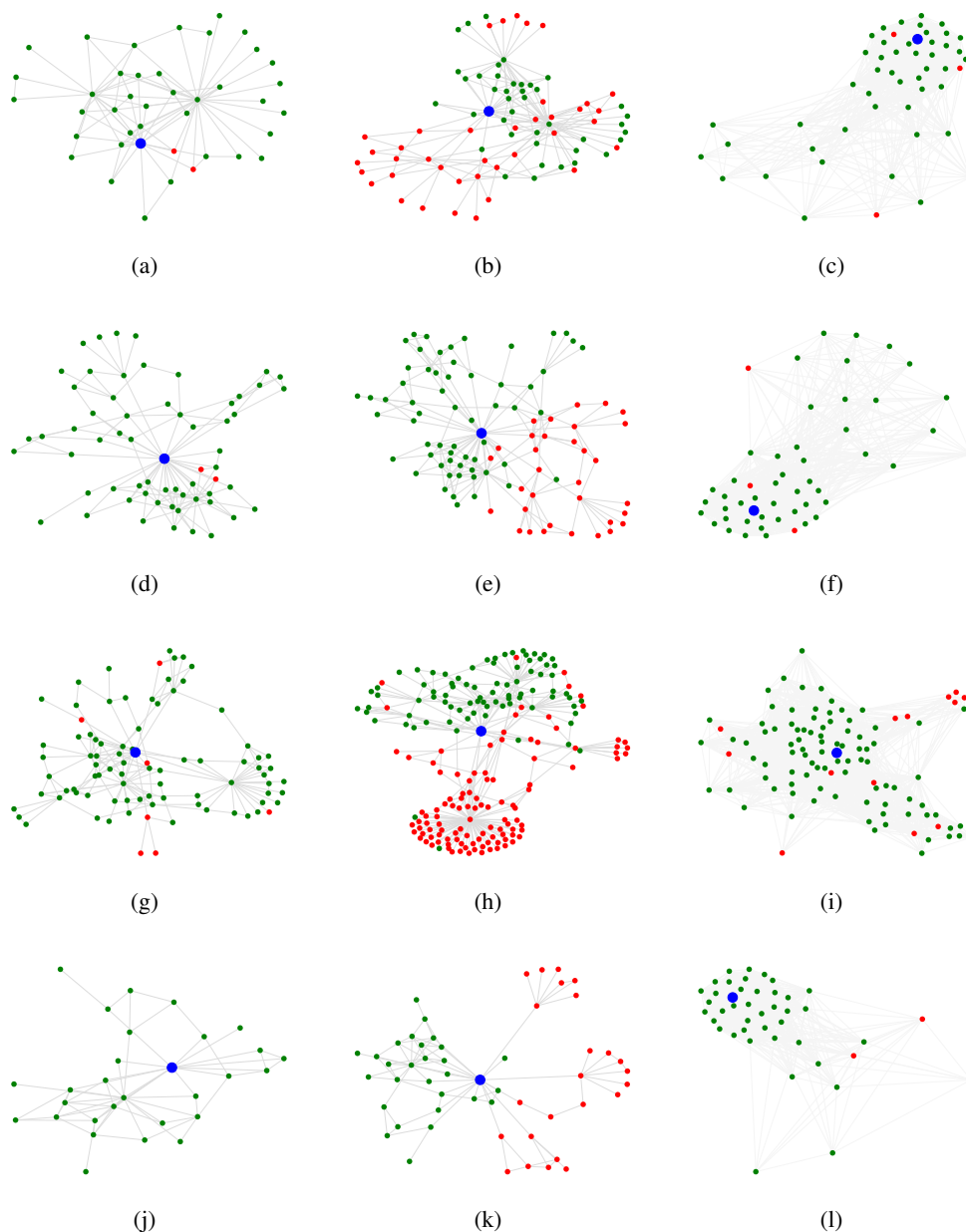
666 **N Additional Visualization Results**

Figure 8: Cora visualizations on a target node (marked in blue) as well as its 1-hop and 2-hop neighbors. Neighbor nodes are marked in green if they have the same label as the target node, and red otherwise. Note that the three graphs in the same row share the same target node (randomly picked), while graphs in different rows focus on different target nodes. **Left:** clean graph. **Middle:** adversarial graph. **Right:** adversarial graph purified by GARNET.

667 We visualize more target nodes and their local structures in Figure 8, which reveals that GARNET
 668 consistently improves the quality of adversarial graph by removing adversarial edges that connect
 669 nodes with different labels. As a result, the adversarial accuracy of backbone GNN models can be
 670 largely improved once they are trained on the GARNET graph.

671 O Homophily Score of GARNET Graph

Table 11: Graph homophily score.

Dataset	Homophilic graphs		Heterophilic graphs	
	Cora	Pubmed	Chameleon	Squirrel
Clean graph	0.80	0.80	0.23	0.22
GARNET graph	0.75	0.72	0.25	0.26

672 We follow Zhu et al. [10] to compute the homophily score per dataset (lower score means more
 673 heterophilic). As shown in Table 11, the GARNET graph is homophilic (heterophilic) if the cor-
 674 responding clean graph is homophilic (heterophilic), which further confirms Theorem 3.3 that our
 675 approach can effectively recover the clean graph structure. As a result, GARNET supports both
 676 homophilic and heterophilic graphs.

677 P Accuracy of Clean Graph Recovery

Table 12: Averaged recall and precision of clean structure recovery over 5 (randomly picked) nodes.

	Recall	Precision
Cora (homophilic graph)	0.94	0.65
Chameleon (heterophilic graph)	0.87	0.59

Apart from visualizing GARNET graph in Figures 5 and 8, we further quantify how well GARNET recovers the clean graph structure. Concretely, given a target node, we first extract nodes within its 2-hop neighbors in the clean graph and GARNET graph (under Metattack with 20% perturbation ratio), respectively. By denoting the extracted nodes by N_{clean} for clean graph and N_{garnet} for GARNET graph, we define the recall score and precision score as follows:

$$Recall = \frac{|N_{clean} \cap N_{garnet}|}{|N_{clean}|}$$

$$Precision = \frac{|N_{clean} \cap N_{garnet}|}{|N_{garnet}|}$$

678 Table 12 shows the averaged recall and precision over 5 nodes on Cora and Chameleon graphs. The
 679 results show that the recall scores are very high for both graphs, which indicates GARNET is able to
 680 accurately recover clean graph structure. The relatively low precision scores indicate that GARNET
 681 also introduces new edges to the graph (i.e., $|N_{garnet}| > |N_{clean}|$). We argue that those new edges
 682 are likely to connect spectrally similar nodes that are far away in the original clean graph, which
 683 enables GARNET to also incorporate global structural information. This explains why GARNET can
 684 sometimes outperform vanilla GNN models on clean heterophilic graphs (shown in Table 2), where
 685 global structural information is very critical for node prediction.

686 Q Further Discussion on Graph Recovery with PGM

687 Intuitively, if we use more (clean) Laplacian eigenpairs (i.e., a larger r) for constructing the embedding
 688 matrix V based on Definition 3.1, the optimal solution for Equation 2 (i.e., Θ^*) will be closer to
 689 the actual clean graph. In this section, we confirm this intuition based on graph resistance distances
 690 between node pairs. Specifically, consider the following expression for calculating effective-resistance
 691 distances between nodes p and q using all Laplacian eigenvalues/eigenvectors except $\lambda_1 = 0$:

$$\sum_{i=2}^{|V|} \frac{(u_i^T e_{p,q})^2}{\lambda_i}$$

692 Feng [17] has shown that the effective-resistance distance between any node pair on the learned graph
 693 Θ^* (when σ approaches infinity) will fully match the Euclidean distance between the corresponding
 694 data samples (i.e., rows in weighted spectral embedding matrix V in our case). Moreover, it can be
 695 shown that the Euclidean distance between the data samples in our case will match the effective-
 696 resistance distance on the original graph when $r = |\mathcal{V}|$ (with proper normalization on Laplacian
 697 eigenvectors). As a result, the resistance distances on the learned graph Θ^* will be the same as
 698 the ones on the original graph when $r = |\mathcal{V}|$. Moreover, using a larger r value will lead to a more
 699 accurate estimation of the learned (clean) graph.

700 In practice, if r satisfies that $\lambda_r \ll \lambda_{r+1}$, dropping the terms with much larger eigenvalues (i.e., λ_{r+1} ,
 701 $\lambda_{r+2}, \dots, \lambda_{|\mathcal{V}|}$) will not significantly impact the approximation accuracy. A proper r can be effectively
 702 determined based on the strategy proposed in Appendix F. We leave the theoretical guarantee of other
 703 metrics for graph comparison to our future work.

704 **R Connection between kNN Graph and TSVD Graph**

705 Apart from the motivation of constructing a kNN graph as \mathcal{G}_{base} based on Theorem 3.3, we further
 706 motivate the kNN graph construction from the perspective of improving the scalability of TSVD-
 707 based defense methods. Concretely, as previous TSVD-based methods produce a dense (low-rank)
 708 adjacency matrix \hat{A} , they involve dense matrices during GNN training, which has quadratic time/space
 709 complexity and thus cannot scale to large graphs. A potential solution is to sparsify \hat{A} by preserving
 710 the top k largest elements per row. However, naïvely selecting the largest elements of each row
 711 in \hat{A} requires forming/storing \hat{A} first, which still has quadratic time/space complexity. In contrast,
 712 we leverage the (approximate) kNN algorithm to construct the sparsified \hat{A} by taking as input the
 713 weighted spectral embedding V (note that $\hat{A} = VV^T$ based on Proposition 3.2). Consequently, our
 714 kNN graph construction step can also be viewed as a scalable way of sparsifying the dense adjacency
 715 matrix \hat{A} generated by TSVD. Moreover, Theorem 3.3 theoretically guarantees that the sparsified
 716 graph serves as a reasonable \mathcal{G}_{base} for edge pruning.