PROVABLE ROBUSTNESS OF (GRAPH) NEURAL NET WORKS AGAINST DATA POISONING AND BACKDOORS

Anonymous authors

Paper under double-blind review

ABSTRACT

Generalization of machine learning models can be severely compromised by data poisoning, where adversarial changes are applied to the training data. This vulnerability has led to interest in certifying (i.e., proving) that such changes up to a certain magnitude do not affect test predictions. We, for the *first* time, certify Graph Neural Networks (GNNs) against poisoning attacks, including backdoors, targeting the node features of a given graph. Our certificates are white-box and based upon (*i*) the *neural tangent kernel*, which characterizes the training dynamics of sufficiently wide networks; and (*ii*) a novel reformulation of the bilevel optimization problem describing poisoning as a mixed-integer linear program. Consequently, we leverage our framework to provide fundamental insights into the role of graph structure and its connectivity on the worst-case robustness behavior of convolution-based and PageRank-based GNNs. We note that our framework is more general and constitutes the *first* approach to derive white-box poisoning certificates for NNs, which can be of independent interest beyond graph-related tasks.

024 025

026 027

004

010 011

012

013

014

015

016

017

018

019

021

1 INTRODUCTION

Numerous works showcase the vulnerability of modern machine learning models to data poisoning,
where adversarial changes are made to the training data (Biggio et al., 2012; Muñoz-González et al.,
2017; Zügner & Günnemann, 2019a; Wan et al., 2023), as well as backdoor attacks affecting both
training and test sets (Goldblum et al., 2023). Empirical defenses against such threats are continually
at risk of being compromised by future attacks (Koh et al., 2022; Suciu et al., 2018). This motivates
the development of *robustness certificates*, which provide formal guarantees that the prediction for a
given test data point remains unchanged under an assumed perturbation model.

Robustness certificates can be categorized as providing deterministic or probabilistic guarantees, and as being white box, i.e. developed for a particular model, or black box (model-agnostic). While each approach has its strengths and applications (Li et al., 2023), we focus on white-box certificates as 037 they can provide a more direct understanding into the worst-case robustness behavior of commonly used models and architectural choices (Tjeng et al., 2019; Mao et al., 2024; Banerjee et al., 2024). The literature on poisoning certificates is less developed than certifying against test-time (evasion) 040 attacks and we provide an overview and categorization in Table 1. Notably, white-box certificates are 041 currently available only for decision trees (Drews et al., 2020), nearest neighbor algorithms (Jia et al., 042 2022), and naive Bayes classification (Bian et al., 2024). In the case of Neural Networks (NNs), the 043 main challenge in white-box poisoning certification comes from capturing their complex training 044 dynamics. As a result, the current literature reveals that deriving white-box poisoning certificates for NNs, and by extension Graph Neural Networks (GNNs), is still an *unsolved* problem, raising the question if such certificates can at all be practically computed. 046

047In this work, we give a positive answer to this question by developing the first approach towards white-
box certification of NNs against data poisoning and backdoor attacks, and instantiate it for common
convolution-based and PageRank-based GNNs. Concretely, poisoning can be modeled as a bilevel
optimization problem over the training data \mathcal{D} that includes training on \mathcal{D} as its inner subproblem.051To overcome the challenge of capturing the complex training dynamics of NNs, we consider the
Neural Tangent Kernel (NTK) that characterizes the training dynamics of sufficiently wide NNs under
gradient flow (Jacot et al., 2018; Arora et al., 2019). In particular, we leverage the equivalence between
NNs trained using the soft-margin loss and standard soft-margin Support Vector Machines (SVMs)

054 Table 1: Representative selection of data poisoning and backdoor attack certificates. To the best of our knowledge, it contains all white-box works. Our work presents the *first* white-box certificate applicable to (graph) neural networks and Support Vector Machines (SVMs). Poisoning refers to (purely) training-time attacks. A backdoor attack refers to joint training-time and test-time 057 perturbations. Certificates apply to different attack types: (i) Clean-label: modifies the features of the training data; (*ii*) Label-flipping: modifies the labels of the training data; (*iii*) Joint: modifies both features and labels; (iv) General attack: allows (arbitrary) insertion/deletion, i.e., like (iii) but 060 dataset size doesn't need to be constant; (v) Node injection: particular to graph learning, refers to 061 adding nodes with arbitrary features and malicious edges into the graph. It is most related to (iv)062 but does not allow deletion and can't be compared with (i) and (ii). Note that certificates that only 063 certify against (iii) - (v) cannot certify against clean-label or label-flipping attacks individually. 064

| | Determinist | c Certified Models | Pois. | Perturi Backd. | ation Model Attack Type | Applies to Node Cls. | Approach |
|--------------------------|--------------|------------------------|--------------|-------------------|----------------------------|-------------------------|-----------------------------|
| (Ma et al., 2019) | × | Diff. Private Learners | ~ | × | Joint | × | Differential Privacy |
| (Liu et al., 2023) | X | Diff. Private Learners | \checkmark | X | General | × | Differential Privacy |
| (Wang et al., 2020) | Х | Smoothed Classifier | X | \checkmark | Joint | X | Randomized Smoothing |
| (Weber et al., 2023) | × X | Smoothed Classifier | X | \checkmark | Clean-label | × | Randomized Smoothing |
| (Zhang et al., 2022) | a x | Smoothed Classifier | \checkmark | \checkmark | Joint | X | Randomized Smoothing |
| (Lai et al., 2024) | × × | Smoothed Classifier | \checkmark | × | Node Injection | \checkmark | Randomized Smoothing |
| (Jia et al., 2021) | × 3lac | Ensemble Classifier | \checkmark | × | General | × | Ensemble (Majority Vote) |
| (Rosenfeld et al., 2020) | | Smoothed Classifier | ~ | X | Label Flip. | × | Randomized Smoothing |
| (Levine & Feizi, 2021) | \checkmark | Ensemble Classifier | \checkmark | X | Label Flip./General | X | Ensemble (Majority Vote) |
| (Wang et al., 2022) | \checkmark | Ensemble Classifier | \checkmark | X | General | X | Ensemble (Majority Vote) |
| (Rezaei et al., 2023) | \checkmark | Ensemble Classifier | \checkmark | × | General | × | Ensemble (Run-Off Election) |
| (Drews et al., 2020) | × ✓ | Decision Trees | \checkmark | X | General | × | Abstract Interpretation |
| (Meyer et al., 2021) | √ Bo | Decision Trees | \checkmark | X | General | X | Abstract Interpretation |
| (Jia et al., 2022) | ∕ te | k-Nearest Neighbors | \checkmark | X | General | X | Majority Vote |
| (Bian et al., 2024) | ,hi ∧ | Naive Bayes Classifier | \checkmark | X | Clean-label | X | Algorithmic |
| Ours | > \ | NNs & SVMs | \checkmark | \checkmark | Clean-label | \checkmark | NTK & Linear Programming |

076 077

071

073 074 075

078 with the NN's NTKs as kernel matrix (Chen et al., 2021). Using this equivalence, we introduce a 079 novel reformulation of the bilevel optimization problem as a mixed-integer linear program (MILP) that allows to certify test datapoints against poisoning as well as backdoor attacks for sufficiently wide

NNs (see Fig. 1). Although our framework applies to wide NNs in 081 general, solving the MILP scales with the number of labeled training samples. Thus, it is a natural fit for semi-supervised learning tasks, where 083 one can take advantage of the low labeling rate. In this context, we 084 focus on semi-supervised node classification in graphs, where certifying 085 against node feature perturbations is particularly challenging due to the interconnectivity between nodes (Zügner & Günnemann, 2019b; Scholten 087 et al., 2023). Here, our framework provides a general and elegant way 088 to handle this interconnectivity inherent to graph learning, by using the 089 corresponding graph NTKs (Sabanayagam et al., 2023) of various GNNs. our poisoning certificate. 090 Our contributions are:



Figure 1: Illustration of

091 (i) We are the first to certify GNNs in node-classification tasks against poisoning and backdoor 092 attacks targeting node features. Our certification framework called QPCert is introduced in Sec. 3 and leverages the NTK to capture the complex training-dynamics of GNNs. Further, it can be applied to 094 NNs in general and thus, it represents the first approach on *white-box* poisoning certificates for NNs.

095 (ii) Enabled by the white-box nature of our certificate, we conduct the first study into the role of graph 096 data and architectural choices on the worst-case robustness of many widely used GNNs against data poisoning and backdoor attacks (see Sec. 4). We focus on convolution-based and PageRank-based 098 architectures and contribute the derivation of the closed-form NTK for APPNP (Gasteiger et al., 099 2019), GIN (Xu et al., 2018), and GraphSAGE (Hamilton et al., 2017) in App. B. 100

(iii) We contribute a reformulation of the bilevel optimization problem describing poisoning as a 101 MILP when instantiated with kernelized SVMs, allowing for white-box certification of SVMs. While 102 we focus on the NTK as kernel, our strategy can be transferred to arbitrary kernel choices. 103

Notation. We represent matrices and vectors with boldfaced upper and lowercase letters, respectively. 104 v_i and M_{ij} denote *i*-th and *ij*-th entries of v and M, respectively. *i*-th row of M is M_i . I_n is the 105 identity matrix and $\mathbf{1}_{n \times n}$ is the matrix of all 1s of size $n \times n$. We use $\langle ., . \rangle$ for scalar product, $\|.\|_2$ for 106 vector Euclidean norm and matrix Frobenius norm, $\mathbb{1}[.]$ for indicator function, \odot for the Hadamard 107 product, $\mathbb{E}[.]$ for expectation, and [z] for the smallest integer $\geq z$ (ceil). [n] denotes $\{1, 2, \ldots, n\}$.

¹⁰⁸ 2 PRELIMINARIES

110 We are given a partially-labeled graph $\mathcal{G} = (S, X)$ with n nodes and a graph structure matrix 111 $S \in \mathbb{R}_{\geq 0}^{n \times n}$, representing for example, a normalized adjacency matrix. Each node $i \in [n]$ has 113 features $x_i \in \mathbb{R}^d$ of dimension d collected in a node feature matrix $X \in \mathbb{R}^{n \times d}$. We assume labels 114 $y_i \in \{1, \ldots, K\}$ are given for the first $m \leq n$ nodes. Our goal is to perform node classification, 115 either in a transductive setting where the labels of the remaining n - m nodes should be inferred, or 116 nodes is denoted \mathcal{V}_L and the set of unlabeled nodes \mathcal{V}_U .

117 **Perturbation model.** We assume that at training time the adversary \mathcal{A} has control over the features 118 of an ϵ -fraction of nodes and that $[(1 - \epsilon)n]$ nodes are clean. For backdoor attacks, the adversary 119 can also change the features of a test node of interest. Following the semi-verified learning setup 120 introduced in (Charikar et al., 2017), we assume that k < n nodes are known to be uncorrupted. 121 We denote the verified nodes by set \mathcal{V}_V and the nodes that can be potentially corrupted as set \mathcal{U} . 122 We further assume that the strength of A to poison training or modify test nodes is bounded by a 123 budget $\delta \in \mathbb{R}_+$. More formally, \mathcal{A} can choose a perturbed $\tilde{x}_i \in \mathcal{B}_p(x_i) := \{\tilde{x} \mid ||\tilde{x} - x_i||_p \le \delta\}$ for each node *i* under control. We denote the set of all perturbed node feature matrices constructible 124 125 by \mathcal{A} from X as $\mathcal{A}(X)$ and $\mathcal{A}(\mathcal{G}) = \{(S, X) \mid X \in \mathcal{A}(X)\}$. In data poisoning, the *goal* of \mathcal{A} is to maximize misclassification in the test nodes. For backdoor attacks \mathcal{A} aims to induce misclassification 126 only in test nodes that it controls. 127

Learning setup. GNNs are functions f_{θ} with (learnable) parameters $\theta \in \mathbb{R}^{q}$ and L number of layers taking the graph $\mathcal{G} = (S, X)$ as input and outputting a prediction for each node. We consider linear output layers with weights W^{L+1} and denote by $f_{\theta}(\mathcal{G})_{i} \in \mathbb{R}^{K}$ the (unnormalized) logit output associated to node *i*. Note for binary classification $f_{\theta}(\mathcal{G})_{i} \in \mathbb{R}$. We define the architectures such as MLP, GCN (Kipf & Welling, 2017), SGC (Wu et al., 2019), (A)PPNP (Gasteiger et al., 2019) and others in App. A. We focus on binary classes $y_{i} \in \{\pm 1\}$ and refer to App. E for the multi-class case. Following Chen et al. (2021), the parameters θ are learned using the soft-margin loss

135

136

137

152 153 154

$$\mathcal{L}(\theta, \mathcal{G}) = \min_{\theta} \frac{1}{2} \| \mathbf{W}^{(L+1)} \|_2^2 + C \sum_{i=1}^m \max(0, 1 - y_i f_{\theta}(\mathcal{G})_i)$$
(1)

where the second term is the Hinge loss weighted by a regularization $C \in \mathbb{R}_+$. Note that due to its non-differentiability, the NN is trained by subgradient descent. Furthermore, we consider NTK parameterization (Jacot et al., 2018) in which parameters θ are initialized from a standard Gaussian $\mathcal{N}(0, 1/\text{width})$. Under NTK parameterization and sufficiently large width limit, the training dynamics of $f_{\theta}(\mathcal{G})$ are precisely characterized by the NTK defined between nodes *i* and *j* as $Q_{ij} = \mathbf{Q}(\mathbf{x}_i, \mathbf{x}_j) = \mathbb{E}_{\theta}[\langle \nabla_{\theta} f_{\theta}(\mathcal{G})_i, \nabla_{\theta} f_{\theta}(\mathcal{G})_j \rangle] \in \mathbb{R}.$

Equivalence of NN to soft-margin SVM with NTK. Chen et al. (2021) show that training NNs in the infinite-width limit with Eq. (1) is equivalent to training a soft-margin SVM with (sub)gradient descent using the NN's NTK as kernel. Thus, both methods converge to the same solution. We extend this equivalence to GNNs, as detailed in App. C. More formally, let the SVM be defined as $f_{\theta}(\mathcal{G})_i = f_{\theta}^{SVM}(\boldsymbol{x}_i) = \langle \boldsymbol{\beta}, \Phi(\boldsymbol{x}_i) \rangle$ where $\Phi(\cdot)$ is the feature transformation associated to the used kernel and $\theta = \boldsymbol{\beta}$ are the learnable parameters obtained by minimizing $\mathcal{L}(\theta, \mathcal{G})$. Following Chen et al. (2021), we do not include a bias term. To find the optimal $\boldsymbol{\beta}^*$, instead of minimizing Eq. (1) with (sub)gradient descent, we work with the equivalent dual

$$P_1(\boldsymbol{Q}): \min_{\boldsymbol{\alpha}} - \sum_{i=1}^m \alpha_i + \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m y_i y_j \alpha_i \alpha_j Q_{ij} \text{ s.t. } 0 \le \alpha_i \le C \ \forall i \in [m]$$
(2)

with the Lagrange multipliers $\alpha \in \mathbb{R}^m$ and kernel $Q_{ij} = Q(x_i, x_j) \in \mathbb{R}$ computed between all labeled nodes $i \in [m]$. and $j \in [m]$. The optimal dual solution may not be unique and we denote by S(Q) the set of α solving $P_1(Q)$. However, any $\alpha^* \in S(Q)$ corresponds to the same unique $\beta^* = \sum_{i=1}^m y_i \alpha_i^* \Phi(\mathcal{G})_i$ minimizing Eq. (1) (Burges & Crisp, 1999). Thus, the prediction of the SVM for a test node t using the dual is given by $f_{\theta}^{SVM}(\mathbf{x}_t) = \sum_{i=1}^m y_i \alpha_i^* Q_{ti}$ for any $\alpha^* \in S(Q)$, where Q_{ti} is the kernel between a test node t and training node i. By choosing Q to be the NTK of a GNN f_{θ} , the prediction equals $f_{\theta}(\mathcal{G})_t$ if the width of the GNN's hidden layers goes to infinity. Thus, a certificate for the SVM directly translates to a certificate for infinitely-wide GNNs. In the finite-width

QPCERT: OUR CERTIFICATION FRAMEWORK

Poisoning a clean training graph \mathcal{G} can be described as a bilevel problem where an adversary \mathcal{A} tries to find a perturbed $\widetilde{\mathcal{G}} \in \mathcal{A}(\mathcal{G})$ that results in a model θ minimizing an attack objective $\mathcal{L}_{att}(\theta, \widetilde{\mathcal{G}})$:

$$\min_{\widetilde{\mathcal{G}},\theta} \mathcal{L}_{att}(\theta,\widetilde{\mathcal{G}}) \quad \text{s. t.} \quad \widetilde{\mathcal{G}} \in \mathcal{A}(\mathcal{G}) \land \theta \in \underset{\theta'}{\operatorname{arg\,min}} \mathcal{L}(\theta',\widetilde{\mathcal{G}}) \tag{3}$$

Eq. (3) is called an upper-level problem and $\min_{\theta'} \mathcal{L}(\theta', \widetilde{\mathcal{G}})$ the lower-level problem. Now, a sample-wise poisoning certificate can be obtained by solving Eq. (3) with an $\mathcal{L}_{att}(\theta, \mathcal{G})$ chosen to describe if the prediction for a test node t changes compared to the prediction of a model trained on the clean graph. However, this approach is challenging as even the simplest bilevel problems given by a linear lower-level problem embedded in an upper-level linear problem are NP-hard (Jeroslow, 1985). Thus, in this section, we develop a general methodology to reformulate the bilevel (sample-wise) certification problem for kernelized SVMs as a mixed-integer linear program, making certification tractable through the use of highly efficient modern MILP solvers such as Gurobi (Gurobi Optimiza-tion, LLC, 2023) or CPLEX (Cplex, 2009). Our approach can be divided into three steps: (1) The bilevel problem is reduced to a single-level problem by exploiting properties of the quadratic dual $P_1(Q)$; (2) We model $\tilde{\mathcal{G}} \in \mathcal{A}(\mathcal{G})$ by assuming a bound on the effect any $\tilde{\mathcal{G}}$ can have on the elements of the kernel Q. This introduces a relaxation of the bilevel problem from Eq. (3) and allows us to fully express certification as a MILP; (3) In Sec. 3.1, we choose the NTK of different GNNs as kernel and develop bounds on the kernel elements to use in the certificate. In the following, we present our certificate for binary classification where $y_i \in \{\pm 1\} \; \forall i \in [n]$ and transductive learning, where the test node is already part of \mathcal{G} . We generalize it to a multi-class and inductive setting in App. E.

A single-level reformulation. Given an SVM f_{θ}^{SVM} trained on the clean graph \mathcal{G} , its class prediction for a test node t is given by $\operatorname{sgn}(\hat{p}_t) = \operatorname{sgn}(f_{\theta}^{SVM}(\boldsymbol{x}_t))$. If for all $\widetilde{\mathcal{G}} \in \mathcal{A}(\mathcal{G})$ the sign of the prediction does not change if the SVM should be retrained on $\widetilde{\mathcal{G}}$, then we know that the prediction for t is certifiably robust. Thus, the attack objective reads $\mathcal{L}_{att}(\theta, \tilde{\mathcal{G}}) = \operatorname{sgn}(\hat{p}_t) \sum_{i=1}^m y_i \alpha_i \tilde{Q}_{ti}$, where \tilde{Q}_{ti} denotes the kernel computed between nodes t and i on the perturbed graph $\widetilde{\mathcal{G}}$, and indicates robustness if greater than zero. Now, notice that the perturbed graph $\widetilde{\mathcal{G}}$ only enters the training objective Eq. (2) through values of the kernel matrix $\widetilde{Q} \in \mathbb{R}^{n \times n}$. Thus, we introduce the set $\mathcal{A}(Q)$ of all kernel matrices \widetilde{Q} , constructable from $\widetilde{\mathcal{G}} \in \mathcal{A}(\mathcal{G})$. Furthermore, we denote with $\mathcal{S}(\widetilde{Q})$ the optimal solution set to $P_1(\tilde{Q})$. As a result, we can rewrite Eq. (3) for kernelized SVMs as

$$P_{2}(\boldsymbol{Q}): \min_{\boldsymbol{\alpha}, \widetilde{\boldsymbol{Q}}} \operatorname{sgn}(\hat{p}_{t}) \sum_{i=1}^{m} y_{i} \alpha_{i} \widetilde{Q}_{ti} \quad s.t. \quad \widetilde{\boldsymbol{Q}} \in \mathcal{A}(\boldsymbol{Q}) \land \boldsymbol{\alpha} \in \mathcal{S}(\widetilde{\boldsymbol{Q}})$$
(4)

and certify robustness if the optimal solution to $P_2(Q)$ is greater than zero. Crucial in reformulating $P_2(Q)$ into a single-level problem are the Karush–Kuhn–Tucker (KKT) conditions of the lower-level problem $P_1(\tilde{Q})$. Concretely, the KKT conditions of $P_1(\tilde{Q})$ are

$$\forall i \in [m]: \quad \sum_{j=1}^{m} y_i y_j \alpha_j \widetilde{Q}_{ij} - 1 - u_i + v_i = 0 \qquad (Stationarity) \quad (5)$$

$$\alpha_i \ge 0, \ C - \alpha_i \ge 0, \ u_i \ge 0, \ v_i \ge 0$$
 (Primal and Dual feasibility) (6)
$$u_i \alpha_i = 0, \ v_i (C - \alpha_i) = 0$$
 (Complementary slackness) (7)

where $\mathbf{u} \in \mathbb{R}^m$ and $\mathbf{v} \in \mathbb{R}^m$ are Lagrange multipliers. Now, we can state (see App. F for the proof):

Proposition 1. Problem $P_1(\tilde{Q})$ given by Eq. (2) is convex and satisfies strong Slater's constraint. Consequently, the single-level optimization problem $P_3(Q)$ arising from $P_2(Q)$ by replacing $\alpha \in S(\tilde{Q})$ with Eqs. (5) to (7) has the same globally optimal solutions as $P_2(Q)$. A mixed-integer linear reformulation. The computational bottleneck of $P_3(\tilde{Q})$ are the non-linear product terms between continuous variables in the attack objective as well as in Eqs. (5) and (7), making $P_3(\tilde{Q})$ a bilinear problem. Thus, we describe in the following how $P_3(\tilde{Q})$ can be transformed into a MILP. First, the complementary slackness constraints can be linearized by recognizing that they have a combinatorial structure. In particular, $u_i = 0$ if $\alpha_i > 0$ and $v_i = 0$ if $\alpha_i < C$. Thus, introducing binary integer variables s and $t \in \{0, 1\}^m$, we reformulate the constraints in Eq. (7) with big-M constraints as

$$\forall i \in [m] : \quad u_i \le M_{u_i} s_i, \ \alpha_i \le C(1 - s_i), \ s_i \in \{0, 1\},$$

$$v_i \le M_{v_i} t_i, \ C - \alpha_i \le C(1 - t_i), \ t_i \in \{0, 1\}$$

$$(8)$$

223 224 225

where M_{u_i} and M_{v_i} are positive constants. In general, verifying that a certain choice of big-Ms results in a valid (mixed-integer) reformulation of the complementary constraints Eq. (7), i.e., such that no optimal solution to the original bilevel problem is cut off, is at least as hard as solving the bilevel problem itself (Kleinert et al., 2020). This is problematic as heuristic choices can lead to suboptimal solutions to the original problem (Pineda & Morales, 2019). However, additional structure provided by $P_1(\tilde{Q})$ and $P_3(Q)$ together with insights into the optimal solution set allows us to derive valid and small M_{u_i} and M_{v_i} for all $i \in [m]$.

233 Concretely, the adversary \mathcal{A} can only make a bounded change to \mathcal{G} . Thus, the element-wise difference of any $\hat{Q} \in \mathcal{A}(Q)$ to Q will be bounded. As a result, there exist element-wise upper and lower 235 bounds $\widetilde{Q}_{ij}^L \leq \widetilde{Q}_{ij} \leq \widetilde{Q}_{ij}^U$ for all $i, j \in [m] \cup \{t\}$ and valid for any $\widetilde{Q} \in \mathcal{A}(Q)$. In Sec. 3.1 we derive concrete lower and upper bounds for the NTKs corresponding to different common GNNs. This, 236 237 together with $0 \le \alpha_i \le C$, allows us to lower and upper bound $\sum_{j=1}^m y_i y_j \alpha_j \widetilde{Q}_{ij}$ in Eq. (5). Now, 238 given an optimal solution $(\alpha^*, \widetilde{Q}^*, \mathbf{u}^*, \mathbf{v}^*)$ to $P_3(\mathbf{Q})$, observe that either u_i^* or v_i^* are zero, or can be 239 freely varied between any positive values as long as Eq. (5) is satisfied without changing the objective 240 value or any other variable. As a result, one can use the lower and upper bounds on $\sum_{j=1}^{m} y_i y_j \alpha_j \tilde{Q}_{ij}$ to find the minimal value range necessary and sufficient for u_i and v_i , such that Eq. (5) can always be 241 242 satisfied for any α^* and \hat{Q}^* . Consequently, only redundant solutions regarding large u_i^* and v_i^* will 243 be cut off and the optimal solution value stays the same as for $P_3(Q)$, not affecting the certification. 244 The exact M_{u_i} and M_{v_i} depend on the signs of the involved y_i and y_j and are derived in App. G. 245

Now, the remaining non-linearities come from the product terms $\alpha_i \hat{Q}_{ij}$. We approach this by first introducing new variables Z_{ij} for all $i, j \in [m] \cup \{t\}$ and set $Z_{ij} = \alpha_j \tilde{Q}_{ij}$. Then, we replace all product terms $\alpha_j \tilde{Q}_{ij}$ in Eq. (5) and in the objective in Eq. (4) with Z_{ij} . This alone has not changed the fact that the problem is bilinear, only that the bilinear terms have now moved to the definition of Z_{ij} . However, we have access to lower and upper bounds on \tilde{Q}_{ij} . Thus, replacing $Z_{ij} = \alpha_j \tilde{Q}_{ij}$ with linear constraints $Z_{ij} \leq \alpha_j \tilde{Q}_{ij}^U$ and $Z_{ij} \geq \alpha_j \tilde{Q}_{ij}^L$ results in a relaxation to $P_3(Q)$. This resolved all non-linearities and we can write the following theorem.

Theorem 1 (MILP Formulation). Node t is certifiably robust against adversary A if the optimal solution to the following MILP denoted by P(Q) is greater than zero

256 257 258 $\min_{\boldsymbol{\alpha}, \mathbf{u}, \mathbf{v}, \mathbf{s}, \mathbf{t}, \mathbf{Z}} \operatorname{sgn}(\hat{p}_t) \sum_{i=1}^m y_i Z_{ti} \quad s.t.$

254

255

$$\forall i \in [m] : \quad \sum_{j=1}^{m} y_i y_j Z_{ij} - 1 - u_i + v_i = 0, \\ u_i \le M_u s_i, \ \alpha_i \le C(1 - s_i), \ s_i \in \{0, 1\}, \\ \alpha_i \ge 0, \ C - \alpha_i \ge 0, \ u_i \ge 0, \ v_i \ge 0, \\ v_i \le M_v t_i, \ C - \alpha_i \le C(1 - t_i), \ t_i \in \{0, 1\}$$

 $Z_{ij} \le \alpha_j \widetilde{Q}_{ij}^U, \ Z_{ij} \ge \alpha_j \widetilde{Q}_{ij}^L \qquad \forall i \in [m] \cup \{t\}, j \in [m]$

P(Q) includes backdoor attacks through the bounds \hat{Q}_{tj}^L and \hat{Q}_{tj}^U for all $j \in [m]$, which for an adversary \mathcal{A} who can manipulate t will be set different. On computational aspects, P(Q) involves $(m+1)^2 + 5m$ variables out of which 2m are binary. Thus, the number of binary variables, which for a particular problem type mainly defines how long it takes MILP-solvers to solve a problem, scales with the number of labeled samples.

270 3.1 **QPCERT FOR GNNS THROUGH THEIR CORRESPONDING NTKS** 271

272 To certify a specific GNN using our QPCert framework, we need to derive 273 element-wise lower and upper bounds 274 valid for all NTK matrices $Q \in \mathcal{A}(Q)$ 275 of the corresponding network, that are 276 constructable by the adversary. As a first step, we introduce the NTKs for 278 the GNNs of interest before deriving 279 the bounds. While Sabanayagam et al. (2023) provides the NTKs for GCN and 281 SGC with and without skip connections, we derive the NTKs for (A)PPNP, GIN 283 and GraphSAGE in App. B. For clar-284 ity, we present the NTKs for $f_{\theta}(\mathcal{G})$ with hidden layers L = 1 here and the general case for any L in the appendix. For L = 1, the NTKs generalize to the form 287 $Q = \mathbf{M}(\mathbf{\Sigma} \odot \dot{\mathbf{E}})\mathbf{M}^T + \mathbf{M}\mathbf{E}\mathbf{M}^T$ for all 288

Table 2: The NTKs of GNNs have the general form Q = $\mathbf{M}(\boldsymbol{\Sigma} \odot \dot{\mathbf{E}})\mathbf{M}^T + \mathbf{M}\mathbf{E}\mathbf{M}^T$ for L = 1. The definitions of $\mathbf{M}, \boldsymbol{\Sigma}, \mathbf{E}$ and $\dot{\mathbf{E}}$ are given in the table. $\mathbf{Z} = \mathbf{S} + \mathbf{I}_n$ and $\mathbf{T} = ((1+\epsilon)\mathbf{I}_n + \mathbf{A}) \mathbf{X}. \quad \kappa_0(z) = \frac{1}{\pi} (\pi - \arccos(z))$ and $\kappa_1(z) = \frac{1}{\pi} (z (\pi - \arccos(z)) + \sqrt{1 - z^2}).$

| GNN | \mathbf{M} | Σ | E_{ij} | \dot{E}_{ij} |
|-----------|----------------|--|---|--|
| GCN | \mathbf{S} | $\mathbf{S}\mathbf{X}\mathbf{X}^T\mathbf{S}^T$ | $\sqrt{\sum_{ii} \sum_{jj} \kappa_1} \left(\frac{\sum_{ij}}{\sqrt{\sum_{ii} \sum_{jj}}} \right)$ | $\kappa_0\left(\frac{\Sigma_{ij}}{\sqrt{\Sigma_{ii}\Sigma_{jj}}}\right)$ |
| SGC | \mathbf{S} | $\mathbf{S}\mathbf{X}\mathbf{X}^T\mathbf{S}^T$ | Σ_{ij} | 1 |
| GraphSAGE | \mathbf{Z} | $\mathbf{Z}\mathbf{X}\mathbf{X}^T\mathbf{Z}^T$ | $\sqrt{\Sigma_{ii}\Sigma_{jj}}\kappa_1\left(\frac{\Sigma_{ij}}{\sqrt{\Sigma_{ii}\Sigma_{jj}}}\right)$ | $\kappa_0\left(\frac{\Sigma_{ij}}{\sqrt{\Sigma_{ii}\Sigma_{jj}}}\right)$ |
| (A)PPNP | Р | $\mathbf{X}\mathbf{X}^T + 1_{n \times n}$ | $\sqrt{\Sigma_{ii}\Sigma_{jj}}\kappa_1\left(\frac{\Sigma_{ij}}{\sqrt{\Sigma_{ii}\Sigma_{jj}}}\right)$ | $\kappa_0\left(\frac{\Sigma_{ij}}{\sqrt{\Sigma_{ii}\Sigma_{jj}}}\right)$ |
| GIN | \mathbf{I}_n | $\mathbf{T}\mathbf{T}^T + 1_{n\times n}$ | $\sqrt{\Sigma_{ii}\Sigma_{jj}}\kappa_1\left(\frac{\Sigma_{ij}}{\sqrt{\Sigma_{ii}\Sigma_{jj}}}\right)$ | $\kappa_0\left(\frac{\Sigma_{ij}}{\sqrt{\Sigma_{ii}\Sigma_{jj}}}\right)$ |
| MLP | \mathbf{I}_n | $\mathbf{X}\mathbf{X}^T + 1_{n\times n}$ | $\sqrt{\Sigma_{ii}\Sigma_{jj}}\kappa_1\left(\frac{\Sigma_{ij}}{\sqrt{\Sigma_{ii}\Sigma_{jj}}}\right)$ | $\kappa_0\left(\frac{\Sigma_{ij}}{\sqrt{\Sigma_{ii}\Sigma_{jj}}}\right)$ |

289 the networks, with the definitions of $\mathbf{M}, \boldsymbol{\Sigma}, \mathbf{E}$ and \mathbf{E} detailed in Table 2. Thus, it is important to note that the effect of the feature matrix \mathbf{X} , which the adversary can manipulate, enters into the NTK only 290 as a product $\mathbf{X}\mathbf{X}^T$, making this the quantity of interest when bounding the NTK matrix. 291

292 Focusing on $p = \{1, 2, \infty\}$ in the perturbation model $\mathcal{B}_p(\mathbf{x})$ and $\mathbf{\tilde{X}} \in \mathcal{A}(\mathbf{X})$, we first derive the 293 bounds for $\widetilde{\mathbf{X}}\widetilde{\mathbf{X}}^T$ by considering $\mathcal{U} := \{i : i \notin \mathcal{V}_V\}$ to be the set of all unverified nodes that the adversary can potentially control. Particularly, we present the worst-case element-wise lower and 295 upper bounds for $\widetilde{\mathbf{X}}\widetilde{\mathbf{X}}^T = \mathbf{X}\mathbf{X}^T + \Delta$ in terms of Δ in Lemma 1, and Lemmas 2 and 3 in App. D. 296

Lemma 1 (Bounds for Δ , $p = \infty$). Given $\mathcal{B}_{\infty}(\mathbf{x})$ and any $\widetilde{\mathbf{X}} \in \mathcal{A}(\mathbf{X})$, then $\widetilde{\mathbf{X}}\widetilde{\mathbf{X}}^T = \mathbf{X}\mathbf{X}^T + \Delta$ where the worst-case bounds for Δ , $\Delta_{ij}^L \leq \Delta_{ij} \leq \Delta_{ij}^U$ for all i and $j \in [n]$, is 297 298

299 300

301 302

303

304

305

306 307 $\Delta_{ij}^{L} = -\delta \|\mathbf{X}_{j}\|_{1} \mathbb{1}[i \in \mathcal{U}] - \delta \|\mathbf{X}_{i}\|_{1} \mathbb{1}[j \in \mathcal{U}] - \delta^{2} d\mathbb{1}[i \in \mathcal{U} \land j \in \mathcal{U} \land i \neq j]$ $\Delta_{ij}^{U} = \delta \|\mathbf{X}_{j}\|_{1} \mathbb{1}[i \in \mathcal{U}] + \delta \|\mathbf{X}_{i}\|_{1} \mathbb{1}[j \in \mathcal{U}] + \delta^{2} d\mathbb{1}[i \in \mathcal{U} \land j \in \mathcal{U}]$ (9)

The NTK bounds \widetilde{Q}_{ij}^L and \widetilde{Q}_{ij}^U , are now derived by simply propagating the bounds for $\widetilde{\mathbf{X}}\widetilde{\mathbf{X}}^T$ through the NTK formulation since the multipliers and addends are positive. To elaborate, we compute \widetilde{Q}_{ii}^L by substituting $\mathbf{X}\mathbf{X}^T = \mathbf{X}\mathbf{X}^T + \Delta^L$, and likewise for \tilde{Q}_{ij}^U . Only bounding E_{ij} and \dot{E}_{ij} needs special care and our respective approach is discussed in App. D.1. Further, we prove that the bounds are tight in the worst-case in App. D.2. 308

Theorem 2 (NTK bounds are tight). The worst-case NTK bounds are tight for GNNs with linear activations such as SGC and (A)PPNP, and MLP with $\sigma(z) = z$ for $p = \{1, 2, \infty\}$ in $\mathcal{B}_p(\mathbf{x})$.

310 311 312

313

309

EXPERIMENTAL RESULTS 4

314 We present (i) the effectiveness of QPCert in certifying different GNNs using their corresponding 315 NTKs against node feature poisoning and backdoor attacks; (*ii*) insights into the role of graph data in 316 worst-case robustness of GNNs, specifically the importance of graph information and its connectivity; 317 (*iii*) a study of the impact of different architectural components in GNNs on their provable robustness.

318 Dataset. We use the real-world graph dataset Cora-ML (Bojchevski & Günnemann, 2018), where 319 we generate continuous 384-dim. embeddings of the abstracts with a modern sentence transformer¹. 320 Furthermore, for binary classification, we use Cora-ML and another real-world graph WikiCS 321 (Mernyei & Cangea, 2022) and extract the subgraphs defined by the two largest classes. We call 322

¹all-MiniLM-L6-v2 https://huggingface.co/sentence-transformers/ 323 from all-MiniLM-L6-v2



(b) WikiCSb: $p_{adv} = 1$ (c) Cora-MLb: $p_{adv} = 0.1$ (d) WikiCSb: $p_{adv} = 0.1$

331 Figure 2: Poison Labeled (PL) setting for Cora-MLb and WikiCS. (a)-(b): QPCert effectively 332 provides non-trivial guarantees. (a)-(d): All GNNs show higher certified accuracy than an MLP. 333

334 the resulting datasets *Cora-MLb* and *WikiCSb*, respectively. Lastly, we use graphs generated from 335 Contextual Stochastic Block Models (CSBM) (Deshpande et al., 2018) for controlled experiments on graph parameters. We give dataset statistics and information on the random graph generation scheme 336 in H.1. For Cora-MLb and WikiCSb, we choose 10 nodes per class for training, leaving 1215 and 337 4640 unlabeled nodes, respectively. For Cora-ML, we choose 20 training nodes per class resulting in 338 2925 unlabeled nodes. From the CSBM, we sample graphs with 200 nodes and choose 40 per class 339 for training, leaving 120 unlabeled nodes. All results are averaged over 5 seeds (Cora-ML: 3 seeds) 340 and reported with standard deviation. We do not need a separate validation set, as we perform 4-fold 341 cross-validation (CV) for hyperparameter tuning. 342

GNNs and attack. We evaluate GCN, SGC, (A)PPNP, GIN, GraphSAGE, MLP, and the skip 343 connection variants GCN Skip- α and GCN Skip-PC (see App. A). All results concern the infinite-344 width limit and thus, are obtained through training an SVM with the corresponding GNN's NTK 345 and, if applicable, applying QPCert using Gurobi to solve the MILP from Theorem 1. We fix the 346 hidden layers to L = 1, and the results for $L = \{2, 4\}$ are provided in App. I.2. For CSBMs we 347 fix C = 0.01 for comparability between experiments and models in the main section. We find 348 that changing C has little effect on the accuracy but can strongly affect the robustness of different 349 architectures. Other parameters on CSBM and all parameters on real-world datasets are set using 350 4-CV (see App. H.2 for details). The SVM's quadratic dual problem is solved using QPLayer 351 (Bambade et al., 2023), a differentiable quadratic programming solver. Thus, for evaluating tightness 352 regarding graph poisoning, we use APGD (Croce & Hein, 2020) with their reported hyperparameters 353 as attack, but differentiate through the learning process using two different strategies: (i) QPLayer, and (ii) the surrogate model proposed in MetaAttack (Zügner & Günnemann, 2019a). To evaluate 354 backdoor tightness, we use the clean-label backdoor attack from Xing et al. (2024) as well as the 355 above APGD attack, but at test time additionally attacking the target node. 356

357 Adversarial evaluation settings. We categorize four settings of interest. (1) Poison Labeled (PL): 358 The adversary \mathcal{A} can potentially poison the labeled data \mathcal{V}_L . (2) Poison Unlabeled (PU): Especially interesting in a semi-supervised setting is the scenario when \mathcal{A} can poison the unlabeled data \mathcal{V}_U , 359 while the labeled data, usually representing a small curated set of high quality, is known to be 360 clean (Shejwalkar et al., 2023). (3) Backdoor Labeled (BL): Like (1) but the test node is also 361 controlled by \mathcal{A} . (4) Backdoor Unlabeled (BU): Like (2) but again, the test node is controlled 362 by A. Settings (1) and (2) are evaluated transductively, i.e. on the unlabeled nodes \mathcal{V}_U already known at training time. Note that this means for (2) that some test nodes may be corrupted. For the 364 backdoor attack settings (3) and (4) the test node is removed from the graph during training and 365 added inductively at test time. The size of the untrusted potential adversarial node set \mathcal{U} is set in 366 percentage $p_{adv} \in \{0.01, 0.02, 0.05, 0.1, 0.2, 0.5, 1\}$ of the scenario-dependent attackable node set 367 and resampled for each seed. We consider node feature perturbations $\mathcal{B}_{p}(\mathbf{x})$ with $p = \{1, 2, \infty\}$ and 368 provide all results concerning p = 1 in App. I.5 and J.4. In the case of CSBM, δ is set in percentage 369 of 2μ of the underlying distribution, and for real data to absolute values. Our main evaluation metric is *certified accuracy*, referring to the percentage of correctly classified nodes without attack that are 370 provably robust against data poisoning / backdoor attacks of the assumed adversary A. We note 371 that we are the first work to study certificates for clean-label attacks on node features in graphs. 372 In particular, all current black-box certificates do not apply to graph learning or ℓ_p perturbation 373 models (see App. M). Thus, there is no baseline prior work. However, we still compare the certified 374 accuracies presented below with two common poisoning defenses in App. I.7. 375

Non-trivial certificates and On the importance of graph information. We evaluate the effective-376 ness of our certificates in providing non-trivial robustness guarantees. Consider the PL setting where 377 A can poison all labeled nodes ($p_{adv} = 1$) for which a trivial certificate would return 0% certified





Figure 3: Certifiable robustness in the Poisoning Unlabeled (PU) and Backdoor Unlabeled (BU) setting with $p_{adv} = 0.1$ for Cora-MLb and $p_{adv} = 0.2$ for CSBM. We refer to App. L for WikiCSb.



Figure 4: Poison Unlabeled for WikiCSb ($p_{adv} = 0.02$), Cora-MLb ($p_{adv} = 0.1$), CSBM ($p_{adv} = 0.2$). (a) Certified accuracy gain: difference of certified acc. to an MLP. (b) Graph connectivity analysis where $c \times$ is cp and cq in CSBM model. (c) APPNP analysis based on α . (d) Tightness of QPCert.

400 accuracy. Figs. 2a and 2b prove that QPCert returns non-trivial guarantees. Further, they highlight 401 an interesting insight: All GNNs have *significantly better* worst-case robustness behavior than the 402 certified accuracy of an MLP. Thus, leveraging the graph connectivity, significantly improves their 403 certified accuracy, even when faced with perturbations on all labeled nodes. In Figs. 2c and 2d we 404 show that this observation stays consistent for other p_{adv} . Similar results for Cora-ML (App. K) and 405 CSBM (App. I.1) establish that this behavior is *not* dataset-specific.

406 In Fig. 3, we evaluate the poison unlabeled (PU) and backdoor unlabeled (BU) settings for different datasets. When poisoning only unlabeled data (PU), the MLP's training process is not affected by the 407 adversary, as the MLP does not access the unlabeled nodes during training. Thus, this provides a good 408 baseline for our certificate to study GNNs. Again, QPCert provides non-trivial certified robustness 409 beyond the MLP baseline. Close to all GNNs show certified accuracy exceeding the one of an MLP 410 for small to intermediate perturbation budgets ($\delta \leq 0.1$) for Cora-MLb (Fig. 3a) and CSBM (Fig. 3a), 411 with a similar picture for WikiCSb (App. L) and Cora-ML (App. K). Note that the drop in certified 412 accuracy for an MLP stems from the transductive learning setting, in which the MLP is confronted 413 with the potentially perturbed unlabeled training nodes at test time. For WikiCSb, Fig. 4a elucidates in 414 detail the certified accuracy gain of an SGC to an MLP and for other GNNs, see App. L (and App. J.1 415 for Cora-MLb, App. I.1 for CSBM). Concerning backdoor attacks on unlabeled nodes Figs. 3b and 3d 416 show that most GNNs show significantly better certified robustness than an MLP, even so MLP training is not affected by A. We observe similar results for a BL setting for Cora-MLb (App. J.1), 417 WikiCS (App. L), and CSBM (App. I.1). These results show that leveraging graph information 418 can significantly improve certified accuracy across all attack settings. Further, across all evaluation 419 settings and datasets, we find GIN and GraphSAGE to provide the lowest certified accuracies of all 420 GNNs; their most important design difference is choosing a sum-aggregation scheme. We note that a 421 comparison across architecture can be affected by the certificate's tightness and we hypothesize that 422 the high worst-case robustness of SGC compared to other models may be due to the certificate being 423 tighter (Theorem 2). However, this still allows us to derive architectural insights for a specific GNN. 424

On graph connectivity and architectural insights. We exemplify study directions enabled through our certification framework. By leveraging CSBMs, we study the effect of graph connectivity in the poisoning unlabeled setting in Fig. 4b for GCN. Interestingly, we observe an inflection point at perturbation strength $\delta = 0.05$, where higher connectivity leads to higher certified accuracy against small perturbations, whereas higher connectivity significantly worsens certified accuracy for strong perturbations. These trends are consistent across various architectures and attack settings (App. I.2).

431 Secondly, we study the effect of different α choices in APPNP on its certified accuracy in poison unlabeled setting in Fig. 4c. Interestingly, it also shows an inflection point in the perturbation strength

432 $(\delta = 0.1)$, where higher α increases the provable robustness for larger δ , whereas worsens the 433 provable robustness for smaller δ in Cora-MLb. Notably, this phenomenon is unique to the PU 434 setting (see App. J.2) and is similarly observed in CSBM as shown in App. I.2. Although this setup 435 seems to be similar to the connectivity analysis, it is different as the α in APPNP realizes weighted 436 adjacency rather than changing the connectivity of the graph, that is, increasing or decreasing the number of edges in the graph. We compare different normalization choices for S in GCN and SGC in 437 App. J.3. Through these analyses, it is significant to note that our certification framework enables 438 informed architectural choices from the perspective of robustness. 439

440 441

442

5 DISCUSSION AND RELATED WORK

How tight is QPCert? We compute an upper bound on provable robustness using APGD by differentiating through the learning process. The results in Fig. 4d show that the provable robustness bounds are tight for small pertubation budgets δ but less tight for larger δ , demonstrating one limitation (other settings and attacks in App. I.3). While theoretically, the NTK bounds are tight (Theorem 2), the approach of deriving element-wise bounds on Q to model A leading to a relaxation of $P_3(Q)$ can explain the gap between provable robustness and empirical attack. Thus, we are excited about opportunities for future work to improve our approach for modeling A in the MILP P(Q).

450 Is QPCert deterministic or probabilistic? Our certification framework is inherently deterministic, 451 offering deterministic guarantees for kernelized SVMs using the NTK as the kernel. When the width 452 of a NN approaches infinity, QPCert provides a deterministic robustness guarantee for the NN due to 453 the exact equivalence between an SVM with the NN's NTK as kernel and the infinitely wide NN. 454 For sufficiently wide but finite-width NNs this equivalence holds with high probability (Chen et al., 455 2021), making our certificate probabilistic in this context. However, note that this high-probability guarantee is qualitatively different from other methods such as randomized smoothing (Cohen et al., 456 2019), in which the certification approach itself is probabilistic and heavily relies on the number of 457 samplings and thus, inherently introduces randomness. 458

459 Generality of QPCert. While we focus on (G)NNs for graph data, our framework enables white-box 460 poisoning certification of NNs on any data domain. QPCert can be extended to other architecture given the criteria outlined in App. N.4 and other tasks such as graph classification (see App. N.5). 461 Further, it allows for certifying general kernelized SVMs for arbitrary kernel choices if respective 462 kernel bounds as in Sec. 3.1 are derived. To the best of our knowledge, this makes our work 463 the first white-box poisoning certificate for kernelized SVMs. Moreover, the reformulation of the 464 bilevel problem to MILP is directly applicable to any quadratic program that satisfies strong Slater's 465 constraint and certain bounds on the involved variables, hence the name QPCert. Thus extensions to 466 certify quadratic programming layers in NN (Amos & Kolter, 2017) or other quadratic learners are 467 thinkable. Therefore, we believe that our work opens up numerous new avenues of research into the 468 provable robustness against data poisoning. 469

Perturbation model. We study semi-verified learning (Charikar et al., 2017). This is particularly 470 interesting for semi-supervised settings, where often a small fraction of nodes are manually verified 471 and labeled (Shejwalkar et al., 2023), or when learning from the crowd Meister & Valiant (2018); 472 Zeng & Shen (2023). However, this may produce overly pessimistic bounds when large fractions 473 of the training data are unverified, but the adversary can only control a small part of it. We study 474 clean-label attacks bounded by ℓ_p -threat models instead of arbitrary perturbations to nodes controlled 475 by \mathcal{A} . We refer to App. N.1 for a discussion with which commonly studied empirical attacks this 476 threat model aligns. Goldblum et al. (2023) distinctively names studying bounded clean-label attacks 477 as an open problem, as most works assume unrealistically large input perturbations. Exemplary, in Fig. 2a, QPCert allows us to certify robustness against ℓ_p -bounded perturbations applied to all labeled 478 data. Most works on poisoning certification work with so-called 'general attacks' allowing arbitrary 479 modifications of data controlled by the adversary. In the setting studied in Fig. 2a, this would always 480 lead to 0 certified accuracy and being unable to provide non-trivial guarantees. 481

Related work. There is little literature on white-box poisoning certificates (see Table 1), and existing
techniques (Drews et al., 2020; Meyer et al., 2021; Jia et al., 2022; Bian et al., 2024) cannot be
extended to NNs. We summarize the most important related work and refer to App. M for more
details. The bilevel problem Eq. (3) is investigated by several works in the context of developing a
poisoning attack or empirical defense, including for SVMs (Biggio et al., 2012; Xiao et al., 2015; Koh

486 & Liang, 2017; Jagielski et al., 2018). Notably Mei & Zhu (2015) reformulate the bilevel problem 487 $P_2(Q)$ for SVMs to a bilinear single-level problem similar to $P_3(Q)$ but only solve it heuristically for 488 attack generation and do not realize the possibility of a MILP reformulation and certification. There are no poisoning certificates for clean-label attacks against GNNs. (Lai et al., 2024) is the only work 489 490 on poisoning certification of GNNs, but differ incomparably in their threat model and are black-box as well as not applicable to backdoors. Lingam et al. (2024) develops a label poisoning attack for 491 GNNs using the bilevel problem with a regression objective and including NTKs as surrogate models. 492 We note that (Steinhardt et al., 2017) develops statistical bounds on the loss that are not applicable to 493 certify classification. 494

495 **Conclusion.** We derive the first white-box poisoning certificate framework for NNs through their 496 NTKs and demonstrate its effective applicability to semi-supervised node classification tasks common in graph learning. In particular, we show that our certificate generates non-trivial robustness guaran-497 tees and insights into the worst-case poisoning robustness to feature perturbations of a wide range of 498 GNNs. The study on node feature perturbations is of practical concern in many application areas of 499 GNNs such as spam detection (Li et al., 2019) or fake news detection (Hu et al., 2024) (see App. N.3 500 for a more detailed discussion), and certification against them poses unique graph-related challenges 501 due to the interconnectedness of nodes. While we address the robustness to node feature perturbations, 502 certifying against structural perturbations to the graph itself remains an open, complex, but important problem and we refer to App. N.2 for a technical discussion on the arising challenges. Thus, this 504 offers a valuable direction for future research. Furthermore, as is the case with all deterministic 505 certificates (Li et al., 2023), scaling to large datasets remains challenging. Consequently, research on 506 scaling deterministic certificates is an impactful avenue for future work. 507

6 ETHICS STATEMENT

Our method represents a robustness certificate for white-box models. This allows a more informed decision when it comes to the safety aspects of currently used models. However, insights into worst-case robustness can be used for good but potentially also by malicious actors. We strongly believe that research about the limitations of existing models is crucial in making models safer and thus, outweighs potential risks. We are not aware of any direct risks coming from our work.

515 516 517

524

526

534

508

509

7 REPRODUCIBILITY STATEMENT

We detail all the experimental setups with the network architectures and hyperparameters in Sec. 4 and app. H.2. The used datasets are open source as mentioned in App. H.1, and the hardware details are discussed in App. H.3. We provide the complete code base with datasets and configuration files to reproduce the experiments in https://figshare.com/s/e155ced9910eb7b3a531. The randomness in the experiments is controlled by setting fixed seeds which are given in the experiment configuration files. The code will be made public upon acceptance.

- References
- Brandon Amos and J. Zico Kolter. Optnet: Differentiable optimization as a layer in neural networks.
 In *International Conference on Machine Learning (ICML)*, 2017.
- Sanjeev Arora, Simon S Du, Wei Hu, Zhiyuan Li, Russ R Salakhutdinov, and Ruosong Wang. On
 exact computation with an infinitely wide neural net. *Advances in Neural Information Processing Systems (NeurIPS)*, 2019.
 - Antoine Bambade, Fabian Schramm, Adrien Taylor, and Justin Carpentier. QPLayer: efficient differentiation of convex quadratic optimization. 2023. URL https://inria.hal.science/ hal-04133055.
- Debangshu Banerjee, Avaljot Singh, and Gagandeep Singh. Interpreting robustness proofs of deep neural networks. In *International Conference on Learning Representations (ICLR)*, 2024.
- 539 Song Bian, Xiating Ouyang, ZHIWEI FAN, and Paris Koutris. Naive bayes classifiers over missing data: Decision and poisoning. In *International Conference on Machine Learning (ICML)*, 2024.

| 540 | |
|-------------------|---|
| 541 542 | Battista Biggio, Blaine Nelson, and Pavel Laskov. Poisoning attacks against support vector machines. In <i>International Conference on Machine Learning (ICML)</i> , 2012. |
| 542 543 544 | Aleksandar Bojchevski and Stephan Günnemann. Deep gaussian embedding of graphs: Unsupervised inductive learning via ranking. In <i>International Conference on Learning Representations</i> , 2018. |
| 545 546 | Christopher J. C. Burges and David Crisp. Uniqueness of the svm solution. In Advances in Neural |
| 5/17 | Information Processing Systems (NeurIPS), 1999. |
| 548 549 | Moses Charikar, Jacob Steinhardt, and Gregory Valiant. Learning from untrusted data. In <i>Annual ACM SIGACT Symposium on Theory of Computing (STOC)</i> , 2017. |
| 550 551 | Ming Chen, Zhewei Wei, Zengfeng Huang, Bolin Ding, and Yaliang Li. Simple and deep graph convolutional networks. In <i>International Conference on Machine Learning (ICML)</i> , 2020. |
| 552 | |
| 553 554 | against data poisoning. In <i>International Conference on Machine Learning</i> , pp. 3299–3319. PMLR, |
| 555 | 2022. |
| 556 557 | Yilan Chen, Wei Huang, Lam Nguyen, and Tsui-Wei Weng. On the equivalence between neural net- work and support vector machine. <i>Advances in Neural Information Processing Systems (NeurIPS)</i> , |
| 558 | 2021. |
| 559 560 | Jeremy Cohen, Elan Rosenfeld, and Zico Kolter. Certified adversarial robustness via randomized smoothing. In <i>International Conference on Machine Learning (ICML)</i> , 2019. |
| 561 | |
| 562 | IBM ILOG Cplex. V12. 1: User's manual for cplex. International Business Machines Corporation, |
| 563 | 2009. |
| 564 | Francesco Croce and Matthias Hein. Reliable evaluation of adversarial robustness with an ensemble |
| 565 | of diverse parameter-free attacks. 2020. |
| 566 | |
| 567 | Enyan Dai, Minhua Lin, Xiang Zhang, and Suhang Wang. Unnoticeable backdoor attacks on graph |
| 568 | neural networks. In <i>Proceedings of the ACM Web Conference 2023</i> , pp. 2263–2273, 2023. |
| 569 | S. Dempe and J. Dutta. Is bilevel programming a special case of a mathematical program |
| 570 | complementarity constraints? <i>Math. Program.</i> , 131:37–48, 2012. doi: https://doi.org/10.1007/ |
| 571 | s10107-010-0342-1. |
| 572 | |
| 573 574 | Yash Deshpande, Subhabrata Sen, Andrea Montanari, and Elchanan Mossel. Contextual stochastic block models. <i>Advances in Neural Information Processing Systems (NeurIPS)</i> , 31, 2018. |
| 575 | Samuel Drews, Aws Albarghouthi, and Loris D'Antoni. Proving data-poisoning robustness in decision |
| 576 577 | trees. In Proceedings of the 41st ACM SIGPLAN Conference on Programming Language Design and Implementation, pp. 1083–1097, 2020. |
| 578 | |
| 579 | Simon S Du, Kangcheng Hou, Russ R Salakhutdinov, Barnabas Poczos, Ruosong Wang, and Keyulu |
| 580 | Xu. Graph neural tangent kernel: Fusing graph neural networks with graph kernels. Advances in |
| 581 | neural information processing systems (NeurIPS), 2019. |
| 582 | Johannes Gasteiger, Aleksandar Bojchevski, and Stephan Günnemann. Predict then propagate: |
| 583 | Graph neural networks meet personalized pagerank. In International Conference on Learning |
| 597 | Representations (ICLR), 2019. |
| 595 | |
| 596 | Jonas Geiping, Liam H Fowl, W. Ronny Huang, Wojciech Czaja, Gavin Taylor, Michael Moeller, |
| 587 | International Conference on Learning Representations (ICLR) 2021 |
| 588 | mernanonal conjerence on Learning Representations (ICLR), 2021. |
| 589 | M. Goldblum, D. Tsipras, C. Xie, X. Chen, A. Schwarzschild, D. Song, A. Madry, B. Li, and |
| 590 | T. Goldstein. Dataset security for machine learning: Data poisoning, backdoor attacks, and |
| 591 | detenses. <i>IEEE Transactions on Pattern Analysis & Machine Intelligence</i> , 45(02):1563–1580, |
| 592 | 2023. ISSN 1939-3539. |
| 593 | Lukas Gosch, Daniel Sturm, Simon Geisler, and Stephan Günnemann. Revisiting robustness in graph machine learning. In <i>International Conference on Learning Representations (ICLR)</i> , 2023. |

| 594 595 596 | Stephan Günnemann. Graph neural networks: Adversarial robustness. In <i>Graph Neural Networks: Foundations, Frontiers, and Applications</i> , pp. 149–176. Springer Singapore, 2022. |
|-------------------|--|
| 597 | Gurobi Optimization, LLC. Gurobi Optimizer Reference Manual, 2023. |
| 598 599 600 | Will Hamilton, Zhitao Ying, and Jure Leskovec. Inductive representation learning on large graphs. Advances in Neural Information Processing Systems (NeurIPS), 2017. |
| 601 602 603 | Christopher Hojny, Shiqiang Zhang, Juan S. Campos, and Ruth Misener. Verifying message-passing neural networks via topology-based bounds tightening. In <i>International Conference on Machine Learning (ICML)</i> , 2024. |
| 604 605 606 | Sanghyun Hong, Nicholas Carlini, and Alexey Kurakin. Diffusion denoising as a certified defense against clean-label poisoning. <i>arXiv preprint arXiv:2403.11981</i> , 2024. |
| 607 608 | Bo Hu, Zhendong Mao, and Yongdong Zhang. An overview of fake news detection: From a new perspective. <i>Fundamental Research</i> , 2024. ISSN 2667-3258. |
| 610 611 612 | W. Ronny Huang, Jonas Geiping, Liam Fowl, Gavin Taylor, and Tom Goldstein. Metapoison: practical general-purpose clean-label data poisoning. In <i>Advances in Neural Information Processing Systems (NeurIPS)</i> , 2020. |
| 613 614 615 | Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and generalization in neural networks. <i>Advances in neural information processing systems</i> , 31, 2018. |
| 616 617 618 | Matthew Jagielski, Alina Oprea, Battista Biggio, Chang Liu, Cristina Nita-Rotaru, and Bo Li. Manipulating machine learning: Poisoning attacks and countermeasures for regression learning. In <i>IEEE Symposium on Security and Privacy (SP)</i> , 2018. |
| 619 620 621 | Robert G Jeroslow. The polynomial hierarchy and a simple model for competitive analysis. <i>Mathematical programming</i> , 1985. |
| 622 623 624 | Jinyuan Jia, Xiaoyu Cao, and Neil Zhenqiang Gong. Intrinsic certified robustness of bagging against data poisoning attacks. In <i>Proceedings of the AAAI conference on artificial intelligence</i> , volume 35, pp. 7961–7969, 2021. |
| 626 627 628 | Jinyuan Jia, Yupei Liu, Xiaoyu Cao, and Neil Zhenqiang Gong. Certified robustness of nearest neighbors against data poisoning and backdoor attacks. In <i>AAAI Conference on Artificial Intelligence (AAAI)</i> , 2022. |
| 629 630 | Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional networks, 2017. |
| 632 633 634 | Thomas Kleinert, Martin Labbé, Fränk Plein, and Martin Schmidt. There's no free lunch: On the hardness of choosing a correct big-m in bilevel optimization. <i>Operations Research</i> , 68 (6): 1716–1721, 2020. doi: https://doi.org/10.1287/opre.2019.1944. |
| 635 636 637 | Pang Wei Koh and Percy Liang. Understanding black-box predictions via influence functions. In <i>International conference on machine learning (ICML)</i> , 2017. |
| 638 639 | Pang Wei Koh, Jacob Steinhardt, and Percy Liang. Stronger data poisoning attacks break data sanitization defenses. <i>Machine Learning</i> , 111(1):1–47, 2022. |
| 640 641 642 | Yuni Lai, Yulin Zhu, Bailin Pan, and Kai Zhou. Node-aware bi-smoothing: Certified robustness against graph injection attacks, 2024. |
| 643 644 645 | Thai Le, Suhang Wang, and Dongwon Lee. Malcom: Generating malicious comments to attack neural fake news detection models. In 2020 IEEE International Conference on Data Mining (ICDM), 2020. |
| 647 | A Levine and S Feizi. Deep partition aggregation: Provable defense against general poisoning attacks. In <i>International Conference on Learning Representations (ICLR)</i> , 2021. |

| 648 649 650 | Ao Li, Zhou Qin, Runshi Liu, Yiqun Yang, and Dong Li. Spam review detection with graph convolutional networks. In <i>International Conference on Information and Knowledge Management</i> (<i>CIKM</i>), 2019. |
|---------------------------------|--|
| 651 652 653 | Linyi Li, Tao Xie, and Bo Li. Sok: Certified robustness for deep neural networks. In <i>IEEE Symposium</i> on Security and Privacy, (SP), 2023. |
| 654 655 656 657 | Vijay Lingam, Mohammad Sadegh Akhondzadeh, and Aleksandar Bojchevski. Rethinking label poisoning for GNNs: Pitfalls and attacks. In <i>International Conference on Learning Representations (ICLR)</i> , 2024. |
| 658 659 660 | Chaoyue Liu, Libin Zhu, and Misha Belkin. On the linearity of large non-linear models: when and why the tangent kernel is constant. <i>Advances in Neural Information Processing Systems (NeurIPS)</i> , 2020. |
| 661 662 663 | Shijie Liu, Andrew C. Cullen, Paul Montague, Sarah M. Erfani, and Benjamin I. P. Rubinstein. Enhancing the antidote: Improved pointwise certifications against poisoning attacks. In AAAI Conference on Artificial Intelligence (AAAI), 2023. |
| 665 666 | Xiaorui Liu, Wei Jin, Yao Ma, Yaxin Li, Liu Hua, Yiqi Wang, Ming Yan, and Jiliang Tang. Elastic graph neural networks. In <i>International Conference on Machine Learning (ICML)</i> , 2021. |
| 667 668 | Yuzhe Ma, Xiaojin Zhu, and Justin Hsu. Data poisoning against differentially-private learners: Attacks and defenses. In <i>International Joint Conference on Artificial Intelligence (IJCAI)</i> , 2019. |
| 670 671 672 | Yuhao Mao, Mark Niklas Mueller, Marc Fischer, and Martin Vechev. Understanding certified training with interval bound propagation. In <i>International Conference on Learning Representations (ICLR)</i> , 2024. |
| 673 674 675 | G.P. McCormick. Computability of global solutions to factorable nonconvex programs: Part i — convex underestimating problems. <i>Mathematical Programming</i> , 10:147–175, 1976. |
| 676 677 | Shike Mei and Xiaojin Zhu. Using machine teaching to identify optimal training-set attacks on machine learners. In AAAI Conference on Artificial Intelligence (AAAI), 2015. |
| 678 679 680 | Michela Meister and Gregory Valiant. A data prism: Semi-verified learning in the small-alpha regime. In <i>Conference On Learning Theory (COLT)</i> , 2018. |
| 681 682 | Péter Mernyei and Cătălina Cangea. Wiki-cs: A wikipedia-based benchmark for graph neural networks. <i>International Conference on Machine Learning (ICML) GRL+ Workshop</i> , 2022. |
| 683 684 685 | Anna Meyer, Aws Albarghouthi, and Loris D'Antoni. Certifying robustness to programmable data bias in decision trees. <i>Advances in Neural Information Processing Systems (NeurIPS)</i> , 2021. |
| 686 687 | Mehryar Mohri, Afshin Rostamizadeh, and Ameet Talwalkar. <i>Foundations of Machine Learning, second edition</i> . The MIT Press, 2018. |
| 688 689 690 691 692 | Luis Muñoz-González, Battista Biggio, Ambra Demontis, Andrea Paudice, Vasin Wongrassamee, Emil C Lupu, and Fabio Roli. Towards poisoning of deep learning algorithms with back-gradient optimization. In <i>Proceedings of the 10th ACM workshop on artificial intelligence and security</i> , pp. 27–38, 2017. |
| 693 694 695 | Salvador Pineda and Juan Miguel Morales. Solving linear bilevel problems using big-ms: Not all that glitters is gold. <i>IEEE Transactions on Power Systems</i> , 34(3):2469–2471, 2019. doi: 10.1109/TPWRS.2019.2892607. |
| 696 697 698 699 | Keivan Rezaei, Kiarash Banihashem, Atoosa Chegini, and Soheil Feizi. Run-off election: Improved provable defense against data poisoning attacks. In <i>International Conference on Machine Learning (ICML)</i> , 2023. |
| 700 701 | Elan Rosenfeld, Ezra Winston, Pradeep Ravikumar, and Zico Kolter. Certified robustness to label- flipping attacks via randomized smoothing. In <i>International Conference on Machine Learning</i> , pp. 8230–8241. PMLR, 2020. |

| 702 703 704 | Mahalakshmi Sabanayagam, Pascal Esser, and Debarghya Ghoshdastidar. Analysis of convolutions, non-linearity and depth in graph neural networks using neural tangent kernel. <i>Transactions on Machine Learning Research (TMLR)</i> , 2023. | |
|--------------------------|---|--|
| 705 706 | Yan Scholten, Jan Schuchardt, Aleksandar Bojchevski, and Stephan Günnemann. Hierarchical | |
| 707 | randomized smoothing. In Advances in Neural Information Processing Systems (NeurIPS), 2023. | |
| 700 | Prithviraj Sen, Galileo Namata, Mustafa Bilgic, Lise Getoor, Brian Galligher, and Tina Eliassi- | |
| 710 | Collective classification in network data. AI Magazine, 29(3):93, Sep. 2008. | |
| 711 | Virat Shejwalkar, Lingjuan Lyu, and Amir Houmansadr. The perils of learning from unlabeled | |
| 712 713 | Backdoor attacks on semi-supervised learning. In International Conference on Computer Vision (ICCV), 2023. | |
| 714 | Amira Saliman and Sarunas Cirdzijauskas Adagraph: Adaptive graph based algorithms for spam | |
| 715 716 | detection in social networks. In <i>Networked Systems</i> , pp. 338–354, 2017. | |
| 717 718 710 | J Michael Steele. <i>The Cauchy-Schwarz master class: an introduction to the art of mathematical inequalities.</i> Cambridge University Press, 2004. | |
| 720 721 | Jacob Steinhardt, Pang Wei W Koh, and Percy S Liang. Certified defenses for data poisoning attacks. Advances in neural information processing systems, 30, 2017. | |
| 722 723 724 725 | Octavian Suciu, Radu Mărginean, Yiğitcan Kaya, Hal Daumé, and Tudor Dumitraş. When does machine learning fail? generalized transferability for evasion and poisoning attacks. In USENIX Conference on Security Symposium (SEC), 2018. | |
| 725 726 727 | Vincent Tjeng, Kai Xiao, and Russ Tedrake. Evaluating robustness of neural networks with mixed integer programming, 2019. | |
| 728 729 730 | Alexander Turner, Dimitris Tsipras, and Aleksander Madry. Clean-label backdoor attacks, 2019. URL https://openreview.net/forum?id=HJg6e2CcK7. | |
| 731 732 | Alexander Wan, Eric Wallace, Sheng Shen, and Dan Klein. Poisoning language models during instruction tuning. In <i>International Conference on Machine Learning (ICML)</i> , 2023. | |
| 733 734 735 | Binghui Wang, Xiaoyu Cao, Jinyuan jia, and Neil Zhenqiang Gong. On certifying robustness against backdoor attacks via randomized smoothing, 2020. | |
| 736 737 738 | Guan Wang, Sihong Xie, Bing Liu, and Philip S. Yu. Identify online store review spammers via social review graph. <i>ACM Trans. Intell. Syst. Technol.</i> , 3(4), 2012. | |
| 739 740 | Haoran Wang, Yingtong Dou, Canyu Chen, Lichao Sun, Philip S. Yu, and Kai Shu. Attacking fake news detectors via manipulating news social engagement. In <i>ACM Web Conference (WWW)</i> , 2023. | |
| 741 742 743 744 | Wenxiao Wang, Alexander Levine, and Soheil Feizi. Improved certified defenses against data poisoning with (deterministic) finite aggregation. In <i>International Conference on Machine Learning (ICML)</i> , 2022. | |
| 745 746 747 | Maurice Weber, Xiaojun Xu, Bojan Karlaš, Ce Zhang, and Bo Li. Rab: Provable robustness against backdoor attacks. In <i>IEEE Symposium on Security and Privacy (SP)</i> , 2023. | |
| 748 749 750 | Felix Wu, Amauri Souza, Tianyi Zhang, Christopher Fifty, Tao Yu, and Kilian Weinberger. Simpli- fying graph convolutional networks. In <i>International Conference on Machine Learning (ICML)</i> , 2019. | |
| 751 752 753 | Zhaohan Xi, Ren Pang, Shouling Ji, and Ting Wang. Graph backdoor. In 30th USENIX Security Symposium (USENIX Security 21), pp. 1523–1540, 2021. | |
| 754 755 | Huang Xiao, Battista Biggio, Gavin Brown, Giorgio Fumera, Claudia Eckert, and Fabio Roli. Is feature selection secure against training data poisoning? In <i>International Conference on Machine Learning (ICML)</i> , pp. 1689–1698. PMLR, 2015. | |

| 756 757 758 759 | Chulin Xie, Yunhui Long, Pin-Yu Chen, Qinbin Li, Sanmi Koyejo, and Bo Li. Unraveling the connections between privacy and certified robustness in federated learning against poisoning attacks. In <i>Proceedings of the 2023 ACM SIGSAC Conference on Computer and Communications Security</i> , 2023. |
|--------------------------|--|
| 760 761 762 | Xiaogang Xing, Ming Xu, Yujing Bai, and Dongdong Yang. A clean-label graph backdoor attack method in node classification task. <i>Knowledge-Based Systems</i> , 304:112433, 2024. |
| 763 764 | Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural networks? In <i>International Conference on Learning Representations (ICLR)</i> , 2018. |
| 765 766 767 | Shiwei Zeng and Jie Shen. Semi-verified pac learning from the crowd. In International Conference on Artificial Intelligence and Statistics (AISTATS), 2023. |
| 768 769 | Xiang Zhang and Marinka Zitnik. Gnnguard: Defending graph neural networks against adversarial attacks. In Advances in Neural Information Processing Systems (NeurIPS), 2020. |
| 770 771 772 773 | Yuhao Zhang, Aws Albarghouthi, and Loris D'Antoni. Bagflip: A certified defense against data poisoning. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho (eds.), Advances in Neural Information Processing Systems (NeurIPS), 2022. |
| 774 775 | Daniel Zügner and Stephan Günnemann. Adversarial attacks on graph neural networks via meta learning. In <i>International Conference on Learning Representations (ICLR)</i> , 2019a. |
| 776 777 778 779 | Daniel Zügner and Stephan Günnemann. Certifiable robustness and robust training for graph convolutional networks. In <i>International Conference on Knowledge Discovery & Data Mining (KDD)</i> , 2019b. |
| 780 781 782 783 | Daniel Zügner and Stephan Günnemann. Certifiable robustness of graph convolutional networks under structure perturbations. In ACM International Conference on Knowledge Discovery & Data Mining (SIGKDD), 2020. |
| 784 | |
| 785 | |
| 786 | |
| 787 | |
| 788 | |
| 789 | |
| 790 | |
| 792 | |
| 793 | |
| 794 | |
| 795 | |
| 796 | |
| 797 | |
| 798 | |
| 799 | |
| 800 | |
| 801 | |
| 802 | |
| 803 | |
| 804 | |
| 805 | |
| 806 | |
| 807 000 | |
| 808 | |

810 А **ARCHITECTURE DEFINITIONS** 811

842

853 854

855

857

812 We consider GNNs as functions f_{θ} with (learnable) parameters $\theta \in \mathbb{R}^q$ and L number of layers taking 813 the graph $\mathcal{G} = (S, X)$ as input and outputs a prediction for each node. We consider linear output 814 layers with weights \mathbf{W}^{L+1} and denote by $f_{\theta}(\mathcal{G})_i \in \mathbb{R}^K$ the (unnormalized) logit output associated 815 to node *i*. In the following, we formally define the (G)NNs such as MLP, GCN (Kipf & Welling, 816 2017), SGC (Wu et al., 2019) and (A)PPNP (Gasteiger et al., 2019) considered in our study. 817

Def. 1 (MLP). The L-layer Multi-Layer Perceptron is defined as $f_{\theta}(\mathcal{G})_i = f_{\theta}^{MLP}(\boldsymbol{x}_i) =$ 818 $W^{L+1}\phi_{\theta}^{(L)}(x_i)$. With $\phi_{\theta}^{l}(x_i) = \sigma(W^{(l)}\phi_{\theta}^{(l-1)}(x_i) + b^{(l)})$ and $\phi_{\theta}^{(0)}(x_i) = x_i$. $W^{(l)} \in \mathbb{R}^{d_{l+1} \times d_l}$ 819 and $\mathbf{b}^{(l)} \in \mathbb{R}^{d_l}$ are the weights/biases of the *l*-th layer with $d_0 = d$ and $d_{L+1} = K$. $\sigma(\cdot)$ is an element-820 821 wise activation function. If not mentioned otherwise, we choose $\sigma(z) = \text{ReLU}(z) = \max\{0, z\}$. 822

Def. 2 (GCN & SGC). A Graph Convolution Network $f_{\theta}^{GCN}(\mathcal{G})$ (Kipf & Welling, 2017) of depth L is defined as $f_{\theta}(\mathcal{G}) = \phi_{\theta}^{(L+1)}(\mathcal{G})$ with $\phi_{\theta}^{(l)}(\mathcal{G}) = \mathbf{S}\sigma(\phi_{\theta}^{(l-1)}(\mathcal{G}))\mathbf{W}^{(l)}$ and $\phi_{\theta}^{(1)}(\mathcal{G}) = \mathbf{S}\mathbf{X}\mathbf{W}^{(1)}$. $\mathbf{W}^{(l)} \in \mathbb{R}^{d_{l-1} \times d_l}$ are the *l*-th layer weights, $d_0 = d$, $d_{L+1} = K$, and $\sigma(z) = ReLU(z)$ applied element-wise. A Simplified Graph Convolution Network $f_{\theta}^{SGC}(\mathcal{G})$ (Wu et al., 2019) is a GCN with 823 824 825 826 linear $\sigma(z) = z$. 827

828 as $f_{\theta}(\mathcal{G}) = \phi_{\theta}^{(L+1)}(\mathcal{G})$ with $\phi_{\theta}^{(l)}(\mathcal{G}) = \sigma(\phi_{\theta}^{(l-1)}(\mathcal{G}))\mathbf{W}_{1}^{(l)} + \mathbf{S}\sigma(\phi_{\theta}^{(l-1)}(\mathcal{G}))\mathbf{W}_{2}^{(l)}$ and $\phi_{\theta}^{(1)}(\mathcal{G}) = \mathbf{X}\mathbf{W}_{1}^{(1)} + \mathbf{S}\mathbf{X}\mathbf{W}_{2}^{(1)}$. $\mathbf{W}_{1}^{(l)}, \mathbf{W}_{2}^{(l)} \in \mathbb{R}^{d_{l-1} \times d_{l}}$ are the l-th layer weights, $d_{0} = d$, $d_{L+1} = K$, and $\sigma(z) = ReLU(z)$ applied element-wise. **S** is fixed to row normalized adjacency (mean aggregator), $\mathbf{D}^{-1}\mathbf{A}$. **Def. 3** (GraphSAGE). The L-layer GraphSAGE $f_{\theta}^{GSAGE}(\mathcal{G})$ (Hamilton et al., 2017) is defined 829 830 831 832 833

Def. 4 (GIN). A one-layer Graph Isomorphism Network $f_{\theta}^{GIN}(\mathcal{G})$ (Xu et al., 2018) is defined as 834 835 $f_{\theta}(\mathcal{G}) = f_{\theta}^{MLP}(\widetilde{\mathbf{G}})$ with $\widetilde{\mathbf{G}} = ((1+\epsilon)\mathbf{I} + \mathbf{A})\mathbf{X}$ where ϵ is a fixed constant and an one-layer ReLU 836 as MLP. 837

Def. 5 ((A)PPNP). The Personalized Propagation of Neural Predictions Network $f_{\theta}^{PPNP}(\mathcal{G})$ 838 Gasteiger et al. (2019) is defined as $f_{\theta}(\mathcal{G}) = \mathbf{P}\mathbf{H}$ where $\mathbf{H}_{i,:} = f_{\theta}^{MLP}(\mathbf{x}_i)$ and $\mathbf{P} = \alpha(\mathbf{I}_n - (1-\alpha)\mathbf{S})^{-1}$. The Approximate PPNP is defined with $\mathbf{P} = (1-\alpha)^K \mathbf{S}^K + \alpha \sum_{i=0}^{K-1} (1-\alpha)^i \mathbf{S}^i$ 839 840 where $\alpha \in [0, 1]$ and $K \in \mathbb{N}$ is a fixed constant. 841

For APPNP and GIN, we consider an MLP with one-layer ReLU activations as given in the default 843 implementation of APPNP. 844

845 Along with the GNNs presented in Definitions 1 to 5, we consider two variants of popular skip connections in GNNs as given a name in Sabanayagam et al. (2023): Skip-PC (pre-convolution), 846 where the skip is added to the features before applying convolution (Kipf & Welling, 2017); and Skip-847 α , which adds the features to each layer without convolving with S (Chen et al., 2020). To facilitate 848 skip connections, we need to enforce constant layer size, that is, $d_i = d_{i-1}$. Therefore, the input 849 layer is transformed using a random matrix W to $\mathbf{H}_0 := \mathbf{X} \mathbf{W}$ of size $n \times h$ where $\mathbf{W}_{ii} \sim \mathcal{N}(0, 1)$ 850 and h is the hidden layer size. Let \mathbf{H}_i be the output of layer i using which we formally define the 851 skip connections as follows. 852

Def. 6 (Skip-PC). In a Skip-PC (pre-convolution) network, the transformed input H_0 is added to the hidden layers before applying the graph convolution **S**, that is, $\forall i \in [L], \phi_{\theta}^{(l)}(\mathcal{G}) =$ $\mathbf{S}\left(\sigma(\phi_{\theta}^{(l-1)}(\mathcal{G})) + \sigma_s(\mathbf{H}_0)\right) \mathbf{W}^{(l)}$, where $\sigma_s(z)$ can be linear or ReLU. 856

Skip-PC definition deviates from Kipf & Welling (2017) because we skip to the input layer instead of 858 the previous one. We define Skip- α as defined in Sabanayagam et al. (2023) similar to Chen et al. 859 (2020).860

861 **Def.** 7 (Skip- α). Given an interpolation coefficient $\alpha \in (0, 1)$, a Skip- α network is defined such that the transformed input \mathbf{H}_0 and the hidden layer are interpolated linearly, that is, $\phi_{\theta}^{(l)}(\mathcal{G}) =$ 863 $\left((1-\alpha)\mathbf{S}\phi_{\theta}^{(l-1)}(\mathcal{G}) + \alpha\sigma_s(\mathbf{H}_0)\right)\mathbf{W}_i \ \forall i \in [L], \text{ where } \sigma_s(z) \text{ can be linear or ReLU.}$

DERIVATION OF NTKS FOR (A)PPNP, GIN AND GRAPHSAGE В

In this section, we derive the NTKs for (A)PPNP, GIN and GraphSAGE, and state the NTKs for GCN and SGC from Sabanayagam et al. (2023).

B.1 NTK FOR (A)PPNP

We derive the closed-form NTK expression for (A)PPNP $f_{\theta}(\mathcal{G})$ (Gasteiger et al., 2019) in this section. The learnable parameters θ are only part of **H**. In practice, $\mathbf{H} = \text{ReLU}(\mathbf{X}\mathbf{W}_1 + \mathbf{B}_1)\mathbf{W}_2 + \mathbf{B}_2$ where node features $\mathbf{X}, \theta = {\mathbf{W}_1 \in \mathbb{R}^{d \times h}, \mathbf{W}_2 \in \mathbb{R}^{h \times K}, \mathbf{B}_1 \in \mathbb{R}^{n \times h}, \mathbf{B}_2 \in \mathbb{R}^{n \times K}}$. Note that in the actual implementation of the MLP, \mathbf{B}_1 is a vector and we consider it to be a matrix by having the same columns so that we can do matrix operations easily. Same for B_2 as well. We give the full architecture with NTK parameterization in the following,

$$f_{\theta}(\mathcal{G}) = \mathbf{P}(\frac{c_{\sigma}}{\sqrt{h}}\sigma(\mathbf{X}\mathbf{W}_1 + \mathbf{B}_1)\mathbf{W}_2 + \mathbf{B}_2)$$
(10)

where $h \to \infty$ and all parameters in θ are initialized as standard Gaussian $\mathcal{N}(0,1)$. c_{σ} is a constant to preserve the input norm (Sabanayagam et al., 2023). We derive for K = 1 as all the outputs are equivalent in expectation. The NTK between nodes i and j is $\underset{\theta \sim \mathcal{N}(0,1)}{\mathbb{E}} [\langle \nabla_{\theta} f_{\theta}(\mathcal{G})_{i}, \nabla_{\theta} f_{\theta}(\mathcal{G})_{j} \rangle].$

Hence, we first write down the gradients for node i following (Arora et al., 2019; Sabanayagam et al., 2023):

$$\begin{aligned} \frac{\partial f_{\theta}(\mathcal{G})_{i}}{\partial \mathbf{W}_{2}} &= \frac{c_{\sigma}}{\sqrt{h}} (\mathbf{P}_{i} \sigma(\mathcal{G}_{1}))^{T} \qquad \qquad ; \mathcal{G}_{1} = \mathbf{X} \mathbf{W}_{1} + \mathbf{B}_{1} \\ \frac{\partial f_{\theta}(\mathcal{G})_{i}}{\partial \mathbf{B}_{2}} &= (\mathbf{P}_{i})^{T} \mathbf{1}_{n} \\ \frac{\partial f_{\theta}(\mathcal{G})_{i}}{\partial \mathbf{W}_{1}} &= \frac{c_{\sigma}}{\sqrt{h}} \mathbf{X}^{T} (\mathbf{P}_{i}^{T} \mathbf{1}_{n} \mathbf{W}_{2}^{T} \odot \dot{\sigma}(\mathcal{G}_{1})) \\ \frac{\partial f_{\theta}(\mathcal{G})_{i}}{\partial \mathbf{B}_{1}} &= \frac{c_{\sigma}}{\sqrt{h}} \mathbf{P}_{i}^{T} \mathbf{1}_{n} \mathbf{W}_{2}^{T} \odot \dot{\sigma}(\mathcal{G}_{1}) \end{aligned}$$

We note that \mathbf{B}_2 has only one learnable parameter for K = 1, but is represented as a vector of size n with all entries the same. Hence, the derivative is simply adding all entries of \mathbf{P}_i . First, we compute the covariance between nodes i and j in \mathcal{G}_1 .

$$\mathbb{E}\left[\left(G_{1}\right)_{ik}\left(G_{1}\right)_{jk'}\right] = \mathbb{E}\left[\left(\mathbf{X}\mathbf{W}_{1} + \mathbf{B}_{1}\right)_{ik}\left(\mathbf{X}\mathbf{W}_{1} + \mathbf{B}_{1}\right)_{jk'}\right]$$

Since the expectation is over \mathbf{W}_1 and \mathbf{B}_1 and all entries are $\sim \mathcal{N}(0,1)$, and i.i.d, the cross terms will be 0 in expectation. Also, for $k \neq k'$, it is 0. Therefore, it gets simplified to

$$\mathbb{E}\left[\left(G_{1}\right)_{ik}\left(G_{1}\right)_{jk}\right] = \mathbb{E}\left[\mathbf{X}_{i}\mathbf{W}_{1}\mathbf{W}_{1}^{T}\mathbf{X}_{j}^{T} + \left(\mathbf{B}_{1}\mathbf{B}_{1}^{T}\right)_{ij}\right]$$
$$= \left(\mathbf{X}\mathbf{X}^{T}\right)_{ij} + 1 = (\Sigma_{1})_{ij} \tag{11}$$

Thus, $\Sigma_1 = \mathbf{X}\mathbf{X}^T + \mathbf{1}_{n \times n}$ and let $(E_1)_{ij} = \mathbb{E}\left[\sigma(\mathcal{G}_1)_i \sigma(\mathcal{G}_1)_j^T\right]$ and $\left(\dot{E}_1\right)_{ij} = \mathbb{E}\left[\dot{\sigma}(\mathcal{G}_1)_i \dot{\sigma}(\mathcal{G}_1)_j^T\right]$ computed using the definitions in Theorem 3 for ReLU non-linearity. Now, we can compute the NTK for each parameter matrix and then sum it up to get the final kernel.

$$\left\langle \frac{\partial f_{\theta}(\mathcal{G})_{i}}{\partial \mathbf{W}_{2}}, \frac{\partial f_{\theta}(\mathcal{G})_{j}}{\partial \mathbf{W}_{2}} \right\rangle = \frac{c_{\sigma}^{2}}{h} \mathbf{P}_{i} \sigma(\mathcal{G}_{1}) \sigma(\mathcal{G}_{1})^{T} \mathbf{P}_{j}^{T}$$

$$\stackrel{h \to \infty}{=} c_{\sigma}^{2} \mathbf{P}_{i} \mathbb{E} \left[\sigma(\mathcal{G}_{1}) \sigma(\mathcal{G}_{1})^{T} \right] \mathbf{P}_{j}^{T} = c_{\sigma}^{2} \mathbf{P}_{i} \mathbf{E}_{1} \mathbf{P}_{j}^{T}$$
(12)

916
917
$$\left\langle \frac{\partial f_{\theta}(\mathcal{G})_{i}}{\partial \mathbf{B}_{2}}, \frac{\partial f_{\theta}(\mathcal{G})_{j}}{\partial \mathbf{B}_{2}} \right\rangle = \mathbf{P}_{i} \mathbf{1}_{n \times n} \mathbf{P}_{j}^{T}$$
(13)

$$\left\langle \frac{\partial f_{\theta}(\mathcal{G})_{i}}{\partial \mathbf{B}_{1}}, \frac{\partial f_{\theta}(\mathcal{G})_{j}}{\partial \mathbf{B}_{1}} \right\rangle \stackrel{h \to \infty}{=} c_{\sigma}^{2} \mathbf{P}_{i} (\mathbb{E} \left[\dot{\sigma}(\mathcal{G}_{1}) \dot{\sigma}(\mathcal{G}_{1}) \right]) \mathbf{P}_{j}^{T} = c_{\sigma}^{2} \mathbf{P}_{i} \dot{\mathbf{E}}_{1} \mathbf{P}_{j}^{T}$$
(14)

 $\sum_{b=1}^{n} (\mathbf{X}^T)_{pb} (\mathbf{P}_j^T \mathbf{W}_2^T)_{bq} \dot{\sigma}(\mathcal{G}_1)_{bq} \Big]$

 $\stackrel{h \to \infty}{=} c_{\sigma}^{2} \sum_{a=1,b=1}^{n,n} (\mathbf{X}\mathbf{X}^{T})_{ab} \mathbf{P}_{ia}(\mathbf{P}^{T})_{bj} \mathbb{E} \left[\dot{\sigma}(\mathcal{G}_{1}) \dot{\sigma}(\mathcal{G}_{1}) \right]_{ab}$

 $= c_{\sigma}^{2} \mathbf{P}_{i}(\mathbf{X}\mathbf{X}^{T} \odot \mathbb{E}\left[\dot{\sigma}(\mathcal{G}_{1})\dot{\sigma}(\mathcal{G}_{1})\right])\mathbf{P}_{j}^{T} = c_{\sigma}^{2} \mathbf{P}_{i}(\mathbf{X}\mathbf{X}^{T} \odot \dot{\mathbf{E}})_{1}\mathbf{P}_{j}^{T}$

(15)

$$\left\langle \frac{\partial f_{\theta}(\mathcal{G})_{i}}{\partial \mathbf{W}_{1}}, \frac{\partial f_{\theta}(\mathcal{G})_{j}}{\partial \mathbf{W}_{1}} \right\rangle = \frac{c_{\sigma}^{2}}{h} \sum_{p,q}^{f,h} (\mathbf{X}^{T}(\mathbf{P}_{i}^{T}\mathbf{1}_{n}\mathbf{W}_{2}^{T}\odot\dot{\sigma}(\mathcal{G}_{1})))_{pq} (\mathbf{X}^{T}(\mathbf{P}_{j}^{T}\mathbf{1}_{n}\mathbf{W}_{2}^{T}\odot\dot{\sigma}(\mathcal{G}_{1})))_{pq} \\ = \frac{c_{\sigma}^{2}}{h} \sum_{p=1}^{d} \sum_{q=1}^{h} \Big[\sum_{a=1}^{n} (\mathbf{X}^{T})_{pa} (\mathbf{P}_{i}^{T}\mathbf{W}_{2}^{T})_{aq} \dot{\sigma}(\mathcal{G}_{1})_{aq} \right]$$

Finally, the NTK matrix for the considered (A)PPNP is sum of Eqs. (12) to (15) as shown below.

$$\mathbf{Q} = c_{\sigma}^{2} \left(\mathbf{P} \mathbf{E}_{1} \mathbf{P}^{T} + \mathbf{P} \mathbf{1}_{n \times n} \mathbf{P}^{T} + \mathbf{P} \dot{\mathbf{E}}_{1} \mathbf{P} + \mathbf{P} \left(\mathbf{X} \mathbf{X}^{T} \odot \dot{\mathbf{E}}_{1} \right) \mathbf{P}^{T} \right)$$
$$= c_{\sigma}^{2} \left(\mathbf{P} \left(\mathbf{E}_{1} + \mathbf{1}_{n \times n} \right) \mathbf{P}^{T} + \mathbf{P} \left(\left(\mathbf{X} \mathbf{X}^{T} + \mathbf{1}_{n \times n} \right) \odot \dot{\mathbf{E}}_{1} \right) \mathbf{P}^{T} \right)$$

$$= c_{\sigma}^{2} \left(\mathbf{P} \left(\mathbf{E}_{1} + \mathbf{1}_{n \times n} \right) \mathbf{P}^{T} + \mathbf{P} \left(\mathbf{\Sigma}_{1} \odot \dot{\mathbf{E}}_{1} \right) \mathbf{P}^{T} \right)$$
(16)

Note that c_{σ} is a constant, and it only scales the NTK, so we set it to 1 in our experiments. Since we use a linear output layer without bias term at the end, that is, $\mathbf{B}_2 = 0$, the NTK we use for our experiments is reduced to

$$\mathbf{Q} = \left(\mathbf{P}\mathbf{E}_{1}\mathbf{P}^{T} + \mathbf{P}\left(\mathbf{\Sigma}_{1}\odot\dot{\mathbf{E}}_{1}\right)\mathbf{P}^{T}\right).$$

B.2 NTK FOR GIN

The GIN architecture Definition 4 is similar to APPNP: **P** and **X** in APPNP are Identity and $\tilde{\mathbf{G}}$ in GIN, respectively. Hence the NTK is exactly the same as APPNP with these matrices. Thus, the NTK for GIN is $\mathbf{O} = \mathbf{E}_{i} + \left(\mathbf{\Sigma}_{i} \otimes \dot{\mathbf{E}}_{i} \right)$

$$\mathbf{Q} = \mathbf{E}_{1} + \left(\mathbf{\Sigma}_{1} \odot \mathbf{E}_{1}\right)$$

with $\mathbf{\Sigma}_{1} = \widetilde{\mathbf{G}}\widetilde{\mathbf{G}}^{T} + \mathbf{1}_{n \times n}, \mathbf{E}_{1} = \underset{\mathbf{F} \sim \mathcal{N}(\mathbf{0}, \mathbf{\Sigma}_{1})}{\mathbb{E}} \left[\sigma(\mathbf{F})\sigma(\mathbf{F})^{T}\right] \text{ and } \dot{\mathbf{E}}_{1} = \underset{\mathbf{F} \sim \mathcal{N}(\mathbf{0}, \mathbf{\Sigma}_{1})}{\mathbb{E}} \left[\dot{\sigma}(\mathbf{F})\dot{\sigma}(\mathbf{F})^{T}\right]. \square$

B.3 NTKS FOR GCN AND SGC

We restate the NTK derived in Sabanayagam et al. (2023) for self containment. The GCN of depth Lwith width $d_l \to \infty \ \forall l \in \{1, \dots, L\}$, the network converges to the following kernel when trained with gradient flow.

Theorem 3 (NTK for Vanilla GCN). For the GCN defined in Definition 2, the NTK \mathbf{Q} at depth L and K = 1 is

$$\mathbf{Q}^{(L)} = \sum_{k=1}^{L+1} \underbrace{\mathbf{S}\Big(\dots\mathbf{S}\Big(\mathbf{S}\Big(\mathbf{\Sigma}_k \odot \dot{\mathbf{E}}_k\Big)\mathbf{S}^T \odot \dot{\mathbf{E}}_{k+1}\Big)\mathbf{S}^T \odot \dots \odot \dot{\mathbf{E}}_L\Big)\mathbf{S}^T.$$
(17)

 Here $\Sigma_k \in \mathbb{R}^{n \times n}$ is the co-variance between nodes of layer k, and is given by $\Sigma_1 = \mathbf{S}\mathbf{X}\mathbf{X}^T\mathbf{S}^T$, $\Sigma_k = \mathbf{S}\mathbf{E}_{k-1}\mathbf{S}^T$ with $\mathbf{E}_k = c_{\sigma} \underset{\mathbf{F} \sim \mathcal{N}(\mathbf{0}, \Sigma_k)}{\mathbb{E}} [\sigma(\mathbf{F})\sigma(\mathbf{F})^T]$, $\dot{\mathbf{E}}_k = c_{\sigma} \underset{\mathbf{F} \sim \mathcal{N}(\mathbf{0}, \Sigma_k)}{\mathbb{E}} [\dot{\sigma}(\mathbf{F})\dot{\sigma}(\mathbf{F})^T]$ and $\dot{\mathbf{E}}$

$$\mathbf{E}_{L+1} = \mathbf{1}_{n \times n}$$

972 973

974

977 978 979

985

986

987

988

989

990

991

992

993

$$\begin{pmatrix} E_k \end{pmatrix}_{ij} = \sqrt{(\Sigma_k)_{ii} (\Sigma_k)_{jj}} \kappa_1 \left(\frac{(\Sigma_k)_{ij}}{\sqrt{(\Sigma_k)_{ii} (\Sigma_k)_{jj}}} \right)$$
$$\begin{pmatrix} \dot{E}_k \end{pmatrix}_{ij} = \kappa_0 \left(\frac{(\Sigma_k)_{ij}}{\sqrt{(\Sigma_k)_{ii} (\Sigma_k)_{jj}}} \right),$$

where
$$\kappa_0(z) = \frac{1}{\pi} (\pi - \arccos(z)) \text{ and } \kappa_1(z) = \frac{1}{\pi} (z (\pi - \arccos(z)) + \sqrt{1 - z^2}).$$

B.4 NTK FOR GRAPHSAGE

From the definition of GraphSAGE Definition 3, it is very similar to GCN with row normalization adjacency $S = D^{-1}A$ where A and D are adjacency matrix with self-loop and its corresponding degree matrix. The differences in GraphSAGE are the following: there is no self-loop to the adjacency as $S = D^{-1}A$, and the neighboring node features are weighted differently compared to the node itself using W_1 and W_2 . Given that the NTK is computed as the expectation over weights at initialization and infinite width, both W_1 and W_2 behave similarly. Hence, these weights can be replaced with a single parameter W which transforms the network definition of GraphSAGE to $\phi_{\theta}^{1}(\mathcal{G}) = (\mathbf{I} + \mathbf{S})\mathbf{X}\mathbf{W}^{(1)}$ and similarly $\phi_{\theta}^{l}(\mathcal{G}) = (\mathbf{I} + \mathbf{S})\sigma(\phi_{\theta}^{(l-1)(\mathcal{G})})\mathbf{W}^{(l)}$ with $\mathbf{S} = \mathbf{D}^{-1}\mathbf{A}$. Thus, the NTK for GraphSAGE is the same as GCN with the difference in the graph normalization S. \Box

994 995 996

997

С EQUIVALENCE OF GNNS TO SVMS

998 We show the equivalence between GNNs and SVMs by extending the result from Chen et al. (2021), 999 which showed that an infinite-width NN trained by gradient descent on a soft-margin loss has the 1000 same training dynamics as that of an SVM with the NN's NTK as the kernel. The fulcrum of their proof that directly depends on the NN is that the NTK stays constant throughout the training (refer to 1001 (Chen et al., 2021, Theorem 3.4)). As we consider the same learning setup with only changing the 1002 network to GNNs, it is enough to show that the graph NTKs stay constant throughout training for the 1003 equivalence to hold in this case. 1004

1005 Constancy of Graph NTKs. This constancy of the NTK in the case of infinitely-wide NNs is deeply studied in Liu et al. (2020) and derived the conditions for the constancy as stated in Theorem 4.

1007 **Theorem 4** ((Liu et al., 2020)). The constancy of the NTK throughout the training of the NN holds if 1008 and only if (i) the last layer of the NN is linear; (ii) the Hessian spectral norm $\|\mathbf{H}\|$ of the neural 1009 network with respect to the parameters is small, that is, $\rightarrow 0$ with the width of the network; (iii) the parameters of the network \mathbf{w} during training and at initialization is bounded, that is, parameters at 1010 time t, \mathbf{w}_t , satisfies $\|\mathbf{w}_t - \mathbf{w}_0\|_2 \leq \epsilon$. 1011

1012 Now, we prove the constancy of graph NTKs of the GNNs by showing the three conditions. 1013

1014 (*i*) **Linear last layer.** The GNNs considered in Definitions 1 and 7 have a linear last layer.

1015 (ii) Small Hessian spectral norm. Recollect that we use NTK parameterization for initializing 1016 the network parameters, that is, $\mathcal{N}(0, 1/\text{width})$. This is equivalent to initializing the network with 1017 standard normal $\mathcal{N}(0,1)$ and appropriately normalizing the layer outputs (Arora et al., 2019; Sa-1018 banayagam et al., 2023). To exemplify, the APPNP network definition with the normalization is given in Eq. (10). Similarly for other GNNs, the normalization results in scaling $\phi_{\theta}^{(l)}$ as $\frac{c_{\sigma}}{\sqrt{h}}\phi_{\theta}^{(l)}$ where h is 1019 1020 the width of the layer l. As all our GNNs have a simple matrix multiplication of the graph structure 1021 without any bottleneck layer, the Hessian spectral norm is $\mathcal{O}(\ln h/\sqrt{h})$ as derived for the multilayer 1022 fully connected networks in Liu et al. (2020). Therefore, as $h \to \infty$ the spectral norm $\to 0$. 1023

(*iii*) Bounded parameters. This is dependent only on the optimization of the loss function as derived 1024 in Chen et al. (2021, Lemma D.1). We directly use this result as our loss and the optimization are the 1025 same as (Chen et al., 2021).

With this, we show that the considered GNNs trained by gradient descent on soft-margin loss is equivalent to SVM with the graph NTK as the kernel. \Box

D DERIVATION OF NTK BOUNDS AND THEOREM 2

In this section, we first present the bounds for Δ in the case of p = 2 and p = 1 in $\mathcal{B}_p(\mathbf{x})$ (Lemma 2 and Lemma 3), and then derive Lemmas 1, 2 and 3 and Theorem 2 stated in Sec. 3.1.

Lemma 2 (Bounds for Δ , p = 2). Given $\mathcal{B}_2(\mathbf{x})$ and any $\widetilde{\mathbf{X}} \in \mathcal{A}(\mathbf{X})$, then $\widetilde{\mathbf{X}}\widetilde{\mathbf{X}}^T = \mathbf{X}\mathbf{X}^T + \Delta$ where the worst-case bounds for Δ , $\Delta_{ij}^L \leq \Delta_{ij} \leq \Delta_{ij}^U$ for all i and $j \in [n]$, is

$$\Delta_{ij}^{L} = -\delta \|\mathbf{X}_{j}\|_{2} \mathbb{1}[i \in \mathcal{U}] - \delta \|\mathbf{X}_{i}\|_{2} \mathbb{1}[j \in \mathcal{U}] - \delta^{2} \mathbb{1}[i \in \mathcal{U} \land j \in \mathcal{U} \land i \neq j]$$

$$\Delta_{ij}^{U} = \delta \|\mathbf{X}_{j}\|_{2} \mathbb{1}[i \in \mathcal{U}] + \delta \|\mathbf{X}_{i}\|_{2} \mathbb{1}[j \in \mathcal{U}] + \delta^{2} \mathbb{1}[i \in \mathcal{U} \land j \in \mathcal{U}]$$
(18)

Lemma 3 (Bounds for Δ , p = 1). Given $\mathcal{B}_2(\mathbf{x})$ and any $\mathbf{\tilde{X}} \in \mathcal{A}(\mathbf{X})$, then $\mathbf{\tilde{X}}\mathbf{\tilde{X}}^T = \mathbf{X}\mathbf{X}^T + \Delta$ where the worst-case bounds for Δ , $\Delta_{ij}^L \leq \Delta_{ij} \leq \Delta_{ij}^U$ for all i and $j \in [n]$, is

$$\Delta_{ij}^{L} = -\delta \|\mathbf{X}_{j}\|_{\infty} \mathbb{1}[i \in \mathcal{U}] - \delta \|\mathbf{X}_{i}\|_{\infty} \mathbb{1}[j \in \mathcal{U}] - \delta^{2} \mathbb{1}[i \in \mathcal{U} \land j \in \mathcal{U} \land i \neq j]$$

$$\Delta_{ij}^{U} = \delta \|\mathbf{X}_{j}\|_{\infty} \mathbb{1}[i \in \mathcal{U}] + \delta \|\mathbf{X}_{i}\|_{\infty} \mathbb{1}[j \in \mathcal{U}] + \delta^{2} \mathbb{1}[i \in \mathcal{U} \land j \in \mathcal{U}]$$
(19)

1046 1047

1053

1056

1061

1064 1065

1029

1030 1031

1032

1033

1039

1043

1045

To derive Lemmas 1, 2 and 3, we consider the perturbed feature matrix $\mathbf{\tilde{X}} \in \mathcal{A}(\mathbf{X})$ and derive the worst-case bounds for $\mathbf{\tilde{X}}\mathbf{\tilde{X}}^T$ based on the perturbation model $\mathcal{B}_p(\mathbf{x})$ where $p = \infty$, p = 2 and p = 1in our study. Let's say \mathcal{U} is the set of nodes that are potentially controlled by the adversary $\mathcal{A}(\mathbf{X})$ and $\mathbf{\tilde{X}} = \mathbf{X} + \mathbf{\Gamma} \in \mathbb{R}^{n \times d}$ where $\mathbf{\Gamma}_i$ is the adversarial perturbations added to node *i* by the adversary, therefore, $\|\mathbf{\Gamma}_i\|_p \le \delta$ and $\mathbf{\Gamma}_i > 0$ for $i \in \mathcal{U}$ and $\mathbf{\Gamma}_i = 0$ for $i \notin \mathcal{U}$. Then

$$\begin{aligned} \hat{\mathbf{X}}\hat{\mathbf{X}}^T &= (\mathbf{X} + \mathbf{\Gamma})(\mathbf{X} + \mathbf{\Gamma})^T \\ &= \mathbf{X}\mathbf{X}^T + \mathbf{\Gamma}\mathbf{X}^T + \mathbf{X}\mathbf{\Gamma}^T + \mathbf{\Gamma}\mathbf{\Gamma}^T = \mathbf{X}\mathbf{X}^T + \Delta. \end{aligned}$$
(20)

1057 As a result, it suffices to derive the worst-case bounds for Δ , $\Delta^L \leq \Delta \leq \Delta^U$, for different 1058 perturbations. To do so, our strategy is to bound the scalar products $\langle \Gamma_i, \mathbf{X}_j \rangle$ and $\langle \Gamma_i, \Gamma_j \rangle$ element-1059 wise, hence derive $\Delta_{ij}^L \leq \Delta_{ij} \leq \Delta_{ij}^U$. In the following, we derive Δ_{ij}^L and Δ_{ij}^U for the cases when 1060 $p = \infty, p = 2$ and p = 1 in $\mathcal{B}_p(\mathbf{x})$.

Case (*i*): **Derivation of Lemma 1 for** $p = \infty$. In this case, the perturbation allows $\|\mathbf{X}_i - \mathbf{X}_i\|_{\infty} \le \delta$, then by Hölder's inequality $\langle \mathbf{a}, \mathbf{b} \rangle \le \|\mathbf{a}\|_p \|\mathbf{b}\|_q$ where $\frac{1}{p} + \frac{1}{q} = 1$ for all $p, q \in [1, \infty]$ we have

$$\begin{aligned} |\langle \mathbf{\Gamma}_{i}, \mathbf{X}_{j} \rangle| &\leq \|\mathbf{\Gamma}_{i}\|_{\infty} \|\mathbf{X}_{j}\|_{1} \leq \delta \|\mathbf{X}_{j}\|_{1} \\ |\langle \mathbf{\Gamma}_{i}, \mathbf{\Gamma}_{j} \rangle| &\leq \|\mathbf{\Gamma}_{i}\|_{2} \|\mathbf{\Gamma}_{j}\|_{2} \leq d \|\Delta_{i}\|_{\infty} \|\Delta_{j}\|_{\infty} \leq d\delta^{2}. \end{aligned}$$
(21)

Using Eq. (21), the worst-case lower bound Δ_{ij}^L is the lower bound of $\Gamma \mathbf{X}^T + \mathbf{X}\Gamma^T + \Gamma\Gamma^T$:

$$\Delta_{ij}^{L} = \begin{cases} 0+ & \text{if } i, j \notin \mathcal{U} \\ -\delta \|\mathbf{X}_{j}\|_{1} + & \text{if } i \in \mathcal{U} \\ -\delta \|\mathbf{X}_{i}\|_{1} + & \text{if } j \in \mathcal{U} \\ -\delta^{2}d & \text{if } i, j \in \mathcal{U} \text{ and } i \neq j. \end{cases}$$
(22)

1074

1069 1070 1071

1075 The last case in Eq. (22) is due to the fact that $\langle \Gamma_i, \Gamma_i \rangle \ge 0$, hence $\Delta_{ii}^L = 0$. Finally, the Eq. (22) can be succinctly written using the indicator function as

1077
1078
$$\Delta_{ij}^L = -\delta \|\mathbf{X}_j\|_1 \mathbb{1}[i \in \mathcal{U}] - \delta \|\mathbf{X}_i\|_1 \mathbb{1}[j \in \mathcal{U}] - \delta^2 d\mathbb{1}[i \in \mathcal{U} \land j \in \mathcal{U} \land i \neq j],$$

deriving the lower bound in Lemma 1. Similarly, applying the Hölder's inequality for the worst-case upper bound, we get

1081 1082

108

1084 1085 1086

1095 1096

1124 1125

1127

Thus, we derive Lemma 1 by succinctly writing it as

$$\Delta_{ij}^{U} = \delta \|\mathbf{X}_{j}\|_{1} \mathbb{1}[i \in \mathcal{U}] + \delta \|\mathbf{X}_{i}\|_{1} \mathbb{1}[j \in \mathcal{U}] + \delta^{2} d\mathbb{1}[i \in \mathcal{U} \land j \in \mathcal{U}].$$

(23)

 $\Delta_{ij}^{U} = \begin{cases} 0+ & \text{if } i, j \notin \mathcal{U} \\ \delta \| \mathbf{X}_{j} \|_{1} + & \text{if } i \in \mathcal{U} \\ \delta \| \mathbf{X}_{i} \|_{1} + & \text{if } j \in \mathcal{U} \\ \delta^{2}d & \text{if } i, j \in \mathcal{U}. \end{cases}$

1091 1092 Case (*ii*): Derivation of Lemma 2 for p = 2. The worst-case lower and upper bounds of Δ_{ij} for 1093 p = 2 is derived in the similar fashion as $p = \infty$. Here, the perturbation allows $\|\tilde{\mathbf{X}}_i - \mathbf{X}_i\|_2 \le \delta$. 1094 Hence,

1097
1098
1099

$$\begin{aligned} |\langle \boldsymbol{\Gamma}_i, \mathbf{X}_j \rangle| &\leq \|\boldsymbol{\Gamma}_i\|_2 \|\mathbf{X}_j\|_2 \leq \delta \|\mathbf{X}_j\|_2 \\ |\langle \boldsymbol{\Gamma}_i, \boldsymbol{\Gamma}_j \rangle| &\leq \|\boldsymbol{\Gamma}_i\|_2 \|\boldsymbol{\Gamma}_j\|_2 \leq \delta^2. \end{aligned}$$
(24)

Using Eq. (24), we derive the lower and upper bounds of Δ_{ij} :

$$\Delta_{ij}^{L} = \begin{cases} 0+ & \text{if } i, j \notin \mathcal{U} \\ -\delta ||\mathbf{X}_{j}||_{2} + & \text{if } i \in \mathcal{U} \\ -\delta ||\mathbf{X}_{i}||_{2} + & \text{if } j \in \mathcal{U} \\ -\delta^{2} & \text{if } i, j \in \mathcal{U} \end{cases} \qquad \Delta_{ij}^{U} = \begin{cases} 0+ & \text{if } i, j \notin \mathcal{U} \\ \delta ||\mathbf{X}_{j}||_{2} + & \text{if } i \in \mathcal{U} \\ \delta ||\mathbf{X}_{i}||_{2} + & \text{if } j \in \mathcal{U} \\ \delta^{2} & \text{if } i, j \in \mathcal{U} \end{cases}$$

1108 1109 1109 1110 1110 1110 **Case** (*iii*): **Derivation of Lemma 3 for** p = 1. The worst-case lower and upper bounds of Δ_{ij} for p = 1 is derived in the similar fashion as $p = \infty$. Here, the perturbation allows $\|\tilde{\mathbf{X}}_i - \mathbf{X}_i\|_1 \le \delta$. Hence,

$$\begin{aligned} |\langle \boldsymbol{\Gamma}_{i}, \mathbf{X}_{j} \rangle| &\leq \|\boldsymbol{\Gamma}_{i}\|_{1} \|\mathbf{X}_{j}\|_{\infty} \leq \delta \|\mathbf{X}_{j}\|_{\infty} \\ |\langle \boldsymbol{\Gamma}_{i}, \boldsymbol{\Gamma}_{j} \rangle| &\leq \|\boldsymbol{\Gamma}_{i}\|_{2} \|\boldsymbol{\Gamma}_{j}\|_{2} \leq \|\boldsymbol{\Gamma}_{i}\|_{1} \|\boldsymbol{\Gamma}_{j}\|_{1} \leq \delta^{2}. \end{aligned}$$
(25)

1117 Using Eq. (25), we derive the lower and upper bounds of Δ_{ij} :

$$\Delta_{ij}^{L} = \begin{cases} 0+ & \text{if } i, j \notin \mathcal{U} \\ -\delta ||\mathbf{X}_{j}||_{\infty} + & \text{if } i \in \mathcal{U} \\ -\delta ||\mathbf{X}_{i}||_{\infty} + & \text{if } j \in \mathcal{U} \\ -\delta^{2} & \text{if } i, j \in \mathcal{U} \end{cases} \qquad \Delta_{ij}^{U} = \begin{cases} 0+ & \text{if } i, j \notin \mathcal{U} \\ \delta ||\mathbf{X}_{j}||_{\infty} + & \text{if } i \in \mathcal{U} \\ \delta ||\mathbf{X}_{i}||_{\infty} + & \text{if } j \in \mathcal{U} \\ \delta^{2} & \text{if } i, j \in \mathcal{U} \end{cases}$$

1126 D.1 BOUNDING E_{ij} AND \dot{E}_{ij} IN THE NTK

1128 NTKs for GNNs with non-linear ReLU activation have **E** and **E** with non-linear $\kappa_1(z)$ and $\kappa_0(z)$ 1129 functions in their definitions, respectively. In order to bound the NTK, we need a strategy to bound 1130 these quantities as well. In this section, we discuss our approach to bound E_{ij} and \dot{E}_{ij} through 1131 bounding the functions for any GNN with L layers. For ease of exposition, we ignore the layer 1132 indexing for the terms of interest and it is understood from the context. Recollect that the definitions 1133 of **E** and **E** are based on Σ , which is a linear combination of **S** and the previous layer. So, we consider that at this stage, we already have Σ , Σ^L and Σ^U . Now, we expand the functions in the definition and write E_{ij} and \dot{E}_{ij} using their corresponding Σ as follows:

$$E_{ij} = \frac{\sqrt{\Sigma_{ii}\Sigma_{jj}}}{\pi} \left(\frac{\Sigma_{ij}}{\sqrt{\Sigma_{ii}\Sigma_{jj}}} \left(\pi - \arccos\left(\frac{\Sigma_{ij}}{\sqrt{\Sigma_{ii}\Sigma_{jj}}}\right) \right) + \sqrt{1 - \frac{\Sigma_{ij}^2}{\Sigma_{ii}\Sigma_{jj}}} \right)$$
(26)

$$\dot{E}_{ij} = \frac{1}{\pi} \left(\pi - \arccos\left(\frac{\Sigma_{ij}}{\sqrt{\Sigma_{ii}\Sigma_{jj}}}\right) \right)$$
(27)

1142 We derive the lower and upper bounds for E_{ij} and \dot{E}_{ij} in Algorithm 1.

1145 Algorithm 1 Procedure to compute $E_{ij}^L, E_{ij}^U, \dot{E}_{ij}^L$ and \dot{E}_{ij}^U

$$\begin{array}{l} \mbox{Given } \boldsymbol{\Sigma}, \, \boldsymbol{\Sigma}^L \mbox{ and } \boldsymbol{\Sigma}^U \\ \mbox{Let } s^l = \sqrt{\boldsymbol{\Sigma}_{ii}^L \boldsymbol{\Sigma}_{jj}^L}, \, s^u = \sqrt{\boldsymbol{\Sigma}_{ii}^U \boldsymbol{\Sigma}_{jj}^U} \\ \mbox{if } \boldsymbol{\Sigma}_{ij}^L > 0 \mbox{ then} \\ a^l = \frac{\boldsymbol{\Sigma}_{ij}^L}{s^u}, \, a^u = \frac{\boldsymbol{\Sigma}_{ij}^U}{s^l} \\ \mbox{else} \\ a^l = \frac{\boldsymbol{\Sigma}_{ij}^L}{s^l}, \, a^u = \frac{\boldsymbol{\Sigma}_{ij}^U}{s^u} \\ \mbox{end if} \\ \mbox{if } |\boldsymbol{\Sigma}_{ij}^U| > |\boldsymbol{\Sigma}_{ij}^L| \mbox{ then} \\ b^l = \left(\frac{\boldsymbol{\Sigma}_{ij}^L}{s^u}\right)^2, \, b^u = \left(\frac{\boldsymbol{\Sigma}_{ij}^U}{s^l}\right)^2 \\ \mbox{else} \\ b^l = \left(\frac{\boldsymbol{\Sigma}_{ij}^L}{s^l}\right)^2, \, b^u = \left(\frac{\boldsymbol{\Sigma}_{ij}^U}{s^u}\right)^2 \\ \mbox{end if} \\ E_{ij}^L = \frac{s^l}{\pi} \left(a^l \left(\pi - \arccos\left(a^l\right)\right) + \sqrt{1 - b^u}\right) \\ E_{ij}^U = \frac{s^u}{\pi} \left(\pi - \arccos\left(a^l\right)\right) \\ \mbox{if } L_{ij}^L = \frac{1}{\pi} \left(\pi - \arccos\left(a^u\right)\right) \end{array}$$

D.2 DERIVATION OF THEOREM 2: NTK BOUNDS ARE TIGHT

We analyze the tightness of NTK bounds by deriving conditions on graph $\mathcal{G} = (\mathbf{S}, \mathbf{X})$ when Δ_{ij}^L and Δ_{ij}^U are attainable exactly. As our NTK bounding strategy is based on bounding the adversarial perturbation $\widetilde{\mathbf{X}}\widetilde{\mathbf{X}}^T$ and the non-linear functions $\kappa_0(z)$ and $\kappa_1(z)$, it is easy to see that the bounds with non-linearities cannot be tight. So, we consider only linear GCN (=SGC), (A)PPNP and MLP with linear activations.

1178 Now, we focus on deriving conditions for the given node features X using the classic result on the 1179 equality condition of Hölder's inequality (Steele, 2004), and then analyze the NTK bounds. Steele 1180 (2004, Fig. 9.1) shows that the bounds on $\langle \mathbf{a}, \mathbf{b} \rangle$ using the Höder's inequality is reached when 1181 $|\mathbf{a}_i|^p = |\mathbf{b}_i|^q \frac{\|\mathbf{a}\|_p^p}{\|\mathbf{b}\|_q^q}$. Using this, we analyze

 $\Delta_{ij} = \langle \mathbf{\Gamma}_i, \mathbf{X}_j \rangle + \langle \mathbf{\Gamma}_j, \mathbf{X}_i \rangle + \langle \mathbf{\Gamma}_i, \mathbf{\Gamma}_j \rangle$ (28)

in which we call $\langle \Gamma_i, \Gamma_j \rangle$ as interaction term. Following this analysis, the tightness of NTK bounds is derived below for $p = \infty$ and p = 2.

Case (*i*): $p = \infty$. In this case, the feature bounds in Eq. (21) are tight,

 $\forall j, \ \mathbf{X}_j \neq 0 \text{ and } \forall i, k \ \mathbf{\Gamma}_{ik} = c_i$

where c_i is some constant such that $\|\Gamma_i\|_{\infty} \leq \delta$ so the perturbation budget is satisfied. As a result, the upper bound of Δ_{ij} in Lemma 1 is achieved exactly in the following cases,

(a) <u>Number of adversarial nodes = 1</u>: Here the interaction term in Eq. (28) is 0 for all *i* and *j*. Then for the one adversarial node *i*, there exists $\mathbf{X}_j \in \mathbb{R}^d_+$, one can set $\Gamma_i = +\delta \mathbf{1}_d$ to achieve the upper bound.

(b) <u>Number of adversarial nodes > 1</u>: Here the interaction term is $\neq 0$ for all the adversarial nodes *i* and *j*. Then, for the adversarial nodes *i* and *j* if there exist $\mathbf{X}_i \in \mathbb{R}^d_+$ and $\mathbf{X}_j \in \mathbb{R}^d_+$ then for $\Gamma_i = \Gamma_j = +\delta \mathbf{1}_d$ upper bounds are achieved.

The NTKs with linear activations Q_{ij} achieve the upper bound in these cases. Similarly, the lower bound in Lemma 1 is achieved exactly as discussed in the following,

(a) <u>Number of adversarial nodes = 1</u>: Here the interaction term in Eq. (28) is 0 for all *i* and *j*. Then for the adversarial node *i*, there exists $\mathbf{X}_j \in \mathbb{R}^d_+$, one can set $\Gamma_i = -\delta \mathbf{1}_d$ to achieve the lower bound.

(b) <u>Number of adversarial nodes > 1</u>: Here the interaction term is $\neq 0$ for all the adversarial nodes i and j. Then, for the adversarial nodes i and j if there exist $\mathbf{X}_i \in \mathbb{R}^d_+$ and $\mathbf{X}_j \in \mathbb{R}^d_-$ then for $\Gamma_i = -\delta \mathbf{1}_d$ and $\Gamma_j = +\delta \mathbf{1}_d$,

1205 leading to tight lower bounds of Lemma 1. The lower and upper tight bounds of Δ together lead to 1206 tight NTK bounds for linear activations. Note that there is no need to impose any structural restriction 1207 on the graph S to achieve the tight bounds for NTK.

Case (*ii*):
$$p = 2$$
. In this case, the feature bounds in Eq. (24) are tight.

1210

$$\forall i, j, \mathbf{X}_i$$
 and $\mathbf{\Gamma}_i$ are linearly dependent

and $\|\Gamma_i\|_2 \le \delta$ so the perturbation budget is satisfied. As a result, the upper bound of Δ_{ij} in Lemma 2 is achieved exactly in the following,

(a) <u>Number of adversarial nodes = 1</u>: Here the interaction term in Eq. (28) is 0 for all *i* and *j*. Then for the one adversarial node *i*, and any $\mathbf{X}_j \in \mathbb{R}^d$, one can set $\Gamma_i = +\delta \frac{\mathbf{X}_j}{\|\mathbf{X}_j\|_2}$ to achieve the upper bound.

1217 (b) <u>Number of adversarial nodes > 1</u>: Here the interaction term is $\neq 0$ for all the adversarial nodes 1218 *i* and *j*. Then, for the adversarial nodes *i* and *j*, if there exist $\mathbf{X}_i \in \mathbb{R}^d_+$ and $\mathbf{X}_j \in \mathbb{R}^d_+$ are linearly 1219 dependent, then for $\Gamma_i = +\delta \frac{\mathbf{X}_j}{\|\mathbf{X}_j\|_2}$ and $\Gamma_j = +\delta \frac{\mathbf{X}_i}{\|\mathbf{X}_i\|_2}$ tight upper bound is achieved.

The NTKs with linear activations Q_{ij} achieve the worst-case upper bound in these cases. Similarly, the lower bound in Lemma 2 is achieved exactly as discussed in the following,

(a) <u>Number of adversarial nodes = 1</u>: Here the interaction term in Eq. (28) is 0 for all *i* and *j*. Then for the adversarial node *i*, and any $\mathbf{X}_j \in \mathbb{R}^d$, one can set $\Gamma_i = -\delta \frac{\mathbf{X}_j}{\|\mathbf{X}_j\|_2}$ to achieve the lower bound.

(b) <u>Number of adversarial nodes > 1</u>: Here the interaction term is $\neq 0$ for all the adversarial nodes i and j. Then, for the adversarial nodes i and j, if there exist $\mathbf{X}_i \in \mathbb{R}^d_+$ and $\mathbf{X}_j \in \mathbb{R}^d_-$ are linearly dependent, then for $\Gamma_i = -\delta \frac{\mathbf{X}_j}{\|\mathbf{X}_i\|_2}$ and $\Gamma_i = +\delta \frac{\mathbf{X}_i}{\|\mathbf{X}_i\|_2}$,

leading to tight lower bounds of Lemma 2. The lower and upper tight bounds of Δ together leads to tight NTK bounds for linear activations. Note that there is no need to impose any structural restriction on the graph S to achieve the tight bounds for NTK, same as the $p = \infty$ case. We further note that only one instance of achieving the worst-case bound is stated, and one can construct similar cases, for example by considering opposite signs for the features and perturbations.

Case (*iii*): p = 1. In this case, the feature bounds in Eq. (25) are tight,

 $\forall j, \mathbf{X}_j \text{ and } \forall i, k \mathbf{\Gamma}_{ik} = c \mathbb{1}[k = \arg \max_{k'} \mathbf{X}_j]$

where $c = \delta$ to satisfy $\|\Gamma_i\|_1 \leq \delta$. As a result, the upper bound of Δ_{ij} in Lemma 3 is achieved exactly in the following cases,

(a) <u>Number of adversarial nodes = 1</u>: Here the interaction term in Eq. (28) is 0 for all *i* and *j*. Then for the one adversarial node *i*, for any *j*, $\mathbf{X}_j \neq 0$ and $\arg \max \mathbf{X}_j = k$, one can set $\Gamma_{ik} = \operatorname{sgn}(X_{jk})\delta$ and $\Gamma_{ik'} = 0 \forall k' \neq k$ to achieve the upper bound. 1242 (b) Number of adversarial nodes > 1: Here the interaction term is $\neq 0$ for all the adversarial nodes 1243 *i* and *j*. Then, for the adversarial nodes *i* and *j* if there exist $\arg \max \mathbf{X}_i = \arg \max \mathbf{X}_j = k$ and 1244 $\operatorname{sgn}(X_{ik}) = \operatorname{sgn}(X_{jk})$ then for $\Gamma_{ik} = \Gamma_{jk} = \operatorname{sgn}(X_{ik})\delta$ and $\forall k' \neq k$, $\Gamma_{ik'} = \Gamma_{jk'} = 0$ upper 1245 bounds are achieved.

The NTKs with linear activations Q_{ij} achieve the upper bound in these cases. Similarly, the lower bound in Lemma 3 is achieved exactly as discussed in the following,

(a) <u>Number of adversarial nodes = 1</u>: Here the interaction term in Eq. (28) is 0 for all *i* and *j*. Then for the one adversarial node *i*, for any *j*, $\mathbf{X}_j \neq 0$ and $\arg \max \mathbf{X}_j = k$, one can set $\Gamma_{ik} = -\operatorname{sgn}(X_{jk})\delta$ and $\Gamma_{ik'} = 0 \forall k' \neq k$ to achieve the lower bound.

1252 (b) <u>Number of adversarial nodes > 1</u>: Here the interaction term is $\neq 0$ for all the adversarial nodes *i* and *j*. Then, for the adversarial nodes *i* and *j* if there exist $\arg \max \mathbf{X}_i = \arg \max \mathbf{X}_j = k$ and $\operatorname{sgn}(X_{ik}) = -\operatorname{sgn}(X_{jk})$ then for $\Gamma_{ik} = -\operatorname{sgn}(X_{ik})\delta$, $\Gamma_{jk} = -\operatorname{sgn}(X_{jk})\delta$ and $\forall k' \neq k$, $\Gamma_{ik'} =$ $\Gamma_{jk'} = 0$,

1256leading to tight lower bounds of Lemma 3. The lower and upper tight bounds of Δ together leads to1257tight NTK bounds for linear activations. Again, there is no need to impose any structural restriction1258on the graph S to achieve the tight bounds for NTK.

1259

1261 E MULTI-CLASS CERTIFICATION

1262

1263 In this section, we discuss the certification for multi-class. We abstract the NN and work with 1264 NTK here. Hence, to do multi-class classification using SVM with NTK, we choose One-Vs-All 1265 strategy, where we learn K classifiers. Formally, we learn β^1, \ldots, β^K which has corresponding 1266 duals $\alpha^1, \ldots, \alpha^K$. In order to learn β^c , all samples with class label c are assumed to be positive 1267 and the rest negative. Assume from hence on that for all c, β^c corresponds to the optimal solution 1268 with the corresponding dual α^c . Then the prediction for a node t is $c^* = \arg \max_c \hat{p}_t^c$ where 1269 $\hat{p}_t^c = \sum_{i=1}^m y_i \alpha_i^c Q_{ti}$ where Q is the NTK matrix.

1270 Given this, we propose a simple extension of our binary certification where to certify a node t1271 as provably robust, we minimize the MILP objective in Theorem 1 for the predicted class c^* and 1272 maximize the objective for the remaining K - 1 classes. Finally, certify t to be provably robust only 1273 if the objective for c^* remains maximum. Formally, we state the objective below.

Theorem 5. Node t with original predicted class c^* is certifiably robust against adversary \mathcal{A} if $c' = c^*$ where c' is defined in the following. Using the MILP $P(\mathbf{Q})$ in Theorem 1, we define

$$P(\boldsymbol{Q})_{c} := P(\boldsymbol{Q}) \text{ using } \boldsymbol{\alpha}^{c}, \text{ with the only change in obj. to } (-1)^{\mathbb{1}[c \neq c^{*}]} \sum_{i=1}^{m} y_{i} Z_{ti}$$

$$c_{t}' = \underset{c \in [K]}{\operatorname{arg max}} P(\boldsymbol{Q})_{c}.$$
(29)

1283

1284 1285

1287

1293

1295

1286 F PROOF OF PROPOSITION 1

We restate Proposition 1.

Proposition 1. Problem $P_1(Q)$ given by Eq. (2) is convex and satisfies strong Slater's constraint. Consequently, the single-level optimization problem $P_3(Q)$ arising from $P_2(Q)$ by replacing $\alpha \in S(\tilde{Q})$ with Eqs. (5) to (7) has the same globally optimal solutions as $P_2(Q)$.

1294 Given any $\widetilde{Q} \in \mathcal{A}(Q)$. We prove two lemmas, leading us towards proving Proposition 1.

Lemma 4. Problem $P_1(\tilde{Q})$ is convex.

1296 Proof. The dual problem $P_1^b(\widetilde{Q})$ associated do an SVM with bias term reads

1300

$$P_{1}^{b}(\widetilde{\boldsymbol{Q}}): \min_{\boldsymbol{\alpha}} - \sum_{i=1}^{m} \alpha_{i} + \frac{1}{2} \sum_{i=1}^{m} \sum_{j=1}^{m} y_{i} y_{j} \alpha_{i} \alpha_{j} \widetilde{Q}_{ij} \text{ s.t.} \sum_{i=1}^{m} \alpha_{i} y_{i} = 0 \quad 0 \le \alpha_{i} \le C \quad \forall i \in [m] \quad (30)$$

1301 It is a known textbook result that $P_1^b(\hat{Q})$ is convex and we refer to Mohri et al. (2018) for a proof. 1302 A necessary and sufficient condition for an optimization problem to be convex is that the objective 1303 function as well as all inequality constraints are convex and the equality constraints affine functions. 1304 Furthermore, the domain of the variable over which is optimized must be a convex set. As removing 1305 the bias term of an SVM results in a dual problem $P_1(\tilde{Q})$ which is equivalent to $P_1^b(\tilde{Q})$ only with 1306 the constraint $\sum_{i=1}^m \alpha_i y_i = 0$ removed, the necessary and sufficient conditions for convexity stay 1307 fulfilled.

1308

Now, we define strong Slater's condition for $P_1(Q)$ embedded in the upper-level problem $P_2(Q)$ defined in Eq. (4), which we from here on will call strong Slater's constraint qualification (Dempe & Dutta, 2012).

1312 Def. 8 (Slater's CQ). *The lower-level convex optimization problem* $P_1(\hat{Q})$ *fulfills strong Slater's* **1313** *Constraint Qualification, if for any upper-level feasible* $\tilde{Q} \in \mathcal{A}(Q)$, *there exists a point* $\alpha(\tilde{Q})$ *in the* **1314** *feasible set of* $P_1(\tilde{Q})$ *such that no constraint in* $P_1(\tilde{Q})$ *is active, i.e.* $0 < \alpha(\tilde{Q})_i < C$ for all $i \in [m]$.

Lemma 5. Problem $P_1(\tilde{Q})$ fulfills strong Slater's constraint qualification.

1317

1318 *Proof.* We prove Lemma 5 through a constructive proof. Given any upper-level feasible $\tilde{Q} \in \mathcal{A}(Q)$. 1319 Let α be an optimal solution to $P_1(\tilde{Q})$. We restrict ourselves to cases, where $P_1(\tilde{Q})$ is non-1320 degenerate, i.e. the optimal solution to the SVM f_{θ}^{SVM} corresponds to a weight vector $\beta \neq 0$. Then, 1321 at least for one index $i \in [m]$ it must hold that $\alpha_i > 0$.

Assume that j is the index in [m] with the smallest $\alpha_j > 0$. Let $\epsilon = \alpha_j/m + 1 > 0$. Now, we construct a new α' from α by for each $i \in [m]$ setting:

- If $\alpha_i = 0$, set $\alpha'_i = \epsilon$.
- If $\alpha_i = C$, set $\alpha'_i = C \epsilon$.
- 1327 1328

1324

1325 1326

The new α' fulfills $0 < \alpha'(\widetilde{Q})_i < C$ for all $i \in [m]$. If $P_1(\widetilde{Q})$ is degenerate, set $\alpha'(\widetilde{Q})_i = C/2$ for all $i \in [m]$. This concludes the proof.

1331 (Dempe & Dutta, 2012) establish that any bilevel optimization problem U whose lower-level problem L is convex and fulfills strong Slater's constraint qualification for any upper-level feasible point has the same global solutions as another problem defined by replacing the lower-level problem L in U with L's Karash Kuhn Tucker conditions. This, together with Lemmas 4 and 5 concludes the proof for Proposition 1.

100

1345 1346

G Setting Big-M constraints

Proposition 2 (Big-M's). Replacing the complementary slackness constraints Eq. (7) in $P_3(Q)$ with the big-M constraints given in Eq. (8) does not cut away solution values of $P_3(Q)$, if for any $i \in [m]$, the big-M values fulfill the following conditions. For notational simplicity j: Condition(j) denotes $j \in \{j \in [m] : Condition(j)\}$.

1344 If $y_i = 1$ then

$$M_{u_i} \ge \sum_{j:y_j=1\land \tilde{Q}_{ij}^U \ge 0} C\tilde{Q}_{ij}^U - \sum_{j:y_j=-1\land \tilde{Q}_{ij}^L \le 0} C\tilde{Q}_{ij}^L - 1$$
(31)

1347
$$j:y_j=1 \land Q_{ij}^2 \ge 0 \qquad j:y_j=-1 \land Q_{ij}^2 \ge 0$$
1348
$$M \ge \sum_{i=1}^{N} Q_{ij}^2 \ge 0 \qquad \sum_{i=1}^$$

1348
1349
$$M_{v_i} \ge \sum_{j:y_j = -1 \land \tilde{Q}_{ij}^U \ge 0} C \tilde{Q}_{ij}^U - \sum_{j:y_j = 1 \land \tilde{Q}_{ij}^L \le 0} C \tilde{Q}_{ij}^L + 1$$
(32)

1350 If $y_i = -1$ then

$$M_{u_i} \ge \sum_{j:y_j = -1 \land \widetilde{Q}_{ij}^U \ge 0} C \widetilde{Q}_{ij}^U - \sum_{j:y_j = 1 \land \widetilde{Q}_{ij}^L \le 0} C \widetilde{Q}_{ij}^L - 1$$
(33)

1357

 $M_{v_i} \ge \sum_{j:y_j=1\land \widetilde{Q}_{ij}^U \ge 0} C\widetilde{Q}_{ij}^U - \sum_{j:y_j=-1\land \widetilde{Q}_{ij}^L \le 0} C\widetilde{Q}_{ij}^L + 1$ (34)

1358 To obtain the tightest formulation for P(Q) from the above conditions, we set the big-M's to equal 1359 the conditions. 1360

1361 *Proof.* Denote by UB an upper bound to $\sum_{j=1}^{m} y_i y_j Z_{ij}$ and by LB a lower bound to $\sum_{j=1}^{m} y_i y_j Z_{ij}$. 1362 The existence of these bounds follows from y_i and $y_j \in \{-1, 1\}$ and $Z_{ij} = \alpha_j \tilde{Q}_{ij}$ with $0 \le \alpha_j \le C$ 1363 and $\tilde{Q}_{ij}^L \le \tilde{Q}_{ij} \le \tilde{Q}_{ij}^U$, i.e. the boundedness of all variables.

1365 u_i and v_i need to be able to be set such that $\sum_{j=1}^m y_i y_j Z_{ij} - u_i + v_i = 1$ (see Eq. (5)) can be satisfied given any α^* and \widetilde{Q}^* part of an optimal solution to $P_3(Q)$. By using UB and LB we get the following inequalities:

$$UB - u_i + v_i \ge 1 \tag{35}$$

1371 1372 and

1373

1369 1370

1374 1375 $LB - u_i + v_i \le 1 \tag{36}$

Denote $\sum_{j=1}^{m} y_i y_j Z_{ij}$ by T. Thus, if $T \ge 1$, setting $v_i = 0$ and $u_i \le UB - 1 \land u_i \ge LB - 1$ allows to satisfy Eq. (5). If T < 1, setting $u_i = 0$ and $v_i \le 1 - LB \land v_i \ge 1 - UB$ allows to satisfy Eq. (5). Note that for a given i, we are free to set u_i and v_i to arbitrary positive values, as long as they satisfy Eq. (5), as they don't affect the optimal solution value nor the values of other variables.

Thus, adding $u_i \leq UB - 1$ and $v_i \leq 1 - LB$ as constraints to $P_3(Q)$ does not affect its optimal solution. Consequently, setting $M_{u_i} \geq UB - 1$ and $M_{v_i} \geq 1 - LB$, are valid big-M constraints in the mixed-integer reformulation of the complementary slackness constraints Eq. (7). The UB and LB values depend on the sign of y_i, y_j and the bounds on α_j and \tilde{Q}_{ij} and the right terms in Eqs. (31) to (34) represent the respective UB and LB arising. This concludes the proof.

1385

1386 1387

H ADDITIONAL EXPERIMENTAL DETAILS

1388 H.1 DATASETS 1389

The CSBM implementation is taken from (Gosch et al., 2023) publicly released under MIT license. 1390 Cora-ML taken from (Bojchevski & Günnemann, 2018) is also released under MIT license. Cora-ML 1391 has 2995 nodes with 8158 edges, and 7 classes. It traditionally comes with a 2879 dimensional 1392 discrete bag-of-words node feature embedding from the paper abstract. As we focus on continuous 1393 perturbation models, we use the abstracts provided by (Bojchevski & Günnemann, 2018) together 1394 with all-MiniLM-L6-v2, a modern sentence transformer from https://huggingface.co/ 1395 sentence-transformers/all-MiniLM-L6-v2 to generate 384-dimensional continuous 1396 node-feature embeddings. From Cora-ML, we extract the subgraph defined by the two most largest classes, remove singleton nodes, and call the resulting binary-classification dataset Cora-MLb. It has 1235 nodes and 2601 edges. WikiCSb, created from extracting the two largest classes from WikiCS, 1399 is the largest used dataset with 4660 nodes and 72806 edges.

1400

1401 H.1.1 CSBM SAMPLING SCHEME 1402

1403 A CSBM graph \mathcal{G} with n nodes is iteratively sampled as: (a) Sample label $y_i \sim Bernoulli(1/2) \ \forall i \in [n]$; (b) Sample feature vectors $\mathbf{X}_i | y_i \sim \mathcal{N}(y_i \boldsymbol{\mu}, \sigma^2 \mathbf{I}_d)$; (c) Sample adjacency $A_{ij} \sim Bernoulli(p)$

| 1/0/ | |
|------|--|
| 1404 | if $y_i = y_j$, $A_{ij} \sim Bernoulli(q)$ otherwise, and $A_{ji} = A_{ij}$. Following Gosch et al. (2023) we set |
| 1405 | p, q through the maximum likelihood fit to Cora (Sen et al., 2008) ($p = 3.17\%, q = 0.74\%$), and μ |
| 1406 | element-wise to $K\sigma/2\sqrt{d}$ with $d = \lfloor n/\ln^2(n) \rfloor$, $\sigma = 1$, and $K = 1.5$, resulting in an interesting |
| 1407 | classification scheme where both graph structure and features are necessary for good generalization. |
| 1408 | We sample $n = 200$ and choose 40 nodes per class for training, leaving 120 unlabeled nodes. |
| 1409 | |
| 1410 | H.2 ARCHITECTURES. |
| 1411 | |
| 1412 | We fix S to S_{row} for GUN, SGC, GUN Skip- α and GUN Skip-PC following (Sabanayagam et al., 2022) S for A DDND as its default implementation. From the CNN definitions App. A the graph |
| 1413 | 2025), S_{sym} for AFFIVE as its default implementation. From the ONV definitions App. A, the graph structure for CIN is $(1 + c)\mathbf{I} + \mathbf{A}$ for GraphSAGE is $\mathbf{I} + \mathbf{D}^{-1}\mathbf{A}$. |
| 1415 | Structure for Onv is $(1 + \epsilon)\mathbf{I} + \mathbf{A}$, for OraphiSAOE is $\mathbf{I} + \mathbf{D} - \mathbf{A}$. |
| 1416 | We outline the hyperparameters for Cora-MLb, for CSBM all parameters are mentioned in Sec. 4 |
| 1417 | except the Skip- α for GCN Skip- α which was set to 0.2. |
| 1418 | CON (Barry Name) C 0.75 |
| 1419 | • UCN (Row Norm.): $C = 0.75$ |
| 1420 | • GCN (Sym. Norm.): $C = 1$ |
| 1421 | • SGC (Row Norm.): $C = 0.75$ |
| 1422 | • SGC (Sym Norm.): $C = 0.75$ |
| 1423 | • ADDND (Sym. Norm.): $C = 1$ $\alpha = 0.1$ |
| 1424 | • APPINF (Sym. Norm.). $C = 1, \alpha = 0.1$ |
| 1425 | • MLP: $C = 0.5$ |
| 1426 | • GCN Skip- α : $C = 1$, $\alpha = 0.1$ |
| 1427 | • GCN Skippc: $C = 0.5$ |
| 1420 | |
| 1430 | For Cora-ML, the following hyperparameters were set: |
| 1431 | • CCN (Bow Norme): $C = 0.05$ |
| 1432 | • UCN (Row Norm.): $C = 0.05$ |
| 1433 | • SGC (Row Norm.): $C = 0.0575$ |
| 1434 | • MLP: $C = 0.004$ |
| 1435 | |
| 1436 | For WikiCSb: |
| 1437 | • GCN (Row Norm): $C = 1$ |
| 1430 | $GGG(\mathbf{P}, \mathbf{N}, \mathbf{r}) = 1$ |
| 1440 | • SGC (Row Norm.): $C = 5$ |
| 1441 | • APPNP (Sym. Norm.): $C = 0.75, \alpha = 0.1$ |
| 1442 | • MLP: $C = 0.175$ |
| 1443 | • GCN Skip- α : $C = 0.1$, $\alpha = 0.1$ |
| 1444 | • GCN Skippc: $C = 1$ |
| 1446 | |
| 1447 | Hyperparameters were set using 4-fold cross-validation and and choosing the result with lowest C in the standard deviation of the heat validation accurate to reduce multiple of the NIL Deviation. |
| 1448 | ne standard deviation of the best vandation accuracy, to reduce runtime of the whiLP certification |
| 1449 | process. |
| 1450 | |
| 1451 | 11.5 HARDWARE. |
| 1452 | Experiments are run on CPU using Gurobi on an internal cluster. Experiments for CSBM, Cora-MLb |

Experiments are run on CPU using Gurobi on an internal cluster. Experiments for CSBM, Cora-MLb and WikiCSb do not require more than 15GB of RAM. Cora-ML experiments do not require more than 20GB of RAM. The time to certify a node depends on the size of MILP as well as the structure of the concrete problem. On our hardware, for CSBM and Cora-MLb certifying one node typically takes several seconds up to one minute on a single CPU. For Cora-ML, certifying a node can take between one minute and several hours (≤ 10) using two CPUs depending on the difficulty of the associated MILP.

1458 I ADDITIONAL RESULTS: CSBM

1460

1461

1478

I.1 EVALUATING QPCERT AND IMPORTANCE OF GRAPH INFORMATION

Fig. 5a shows the same result as Fig. 2a from Sec. 4 establishing that including graph information boosts worst-case robustness in CSBM too. This also shows that the result is not dataset-specific. In Fig. 5, we provide the remaining settings in correspondence to Fig. 3, Poison Labeled PL and Backdoor Labeled BL for CSBM. Similarly, the heatmaps showing the certified accuracy gain with respect to MLP is presented in Fig. 6.



Figure 5: Certifiable robustness for different (G)NNs in Poisoning Labeled (PL) and Backdoor Labeled (BL) setting.



1512 I.2 ON GRAPH CONNECTIVITY AND ARCHITECTURAL INSIGHTS

1514 We present the sparsity analysis for SGC and APPNP in (a) and (b) of Fig. 7, showing a similar 1515 observation to GCN in App. I.2. The APPNP α analysis for PU and PL are provided in (c) and (d) 1516 of Fig. 7, showing the inflection point in PU but not in PL. Additionally, we show the influence of 1517 depth, linear vs ReLU, regularization C and row vs symmetric normalized adjacency in Fig. 8.



Figure 7: (a)-(b): Graph connectivity analysis where $c \times$ is cp and cq in CSBM model. GCN is provided in Fig. 4b. (c)-(d): APPNP analysis based on α .



Figure 8: (a): Symmetric and row normalized adjacencies as the choice for S in GCN and SGC. (b): Effect of number of hidden layers L. (c): Linear and relu for the Skip-PC and Skip- α . (d): Regularization C in GCN. All experiments in PU setting and $p_{adv} = 0.2$.

1540 1541 1542

1543

1538

1539

1526

1527

1528

I.3 TIGHTNESS OF QPCERT

1544 First, we present the tightness of QPCert in Figs. 9 and 10 evaluated with our strongest employed 1545 attacks for each setting: For graph poisoning (Fig. 9), APGD is employed with direct differentiation 1546 through the learning process (QPLayer) for the PU setting in Fig. 9b and for the PL setting in Fig. 9a. For the backdoor attack setting (Fig. 10), first a poisoning attack is carried out with APGD (QPLayer) 1547 and then, the respective test node is additionally attacked with APGD in an evasion setting. Fig. 10b 1548 shows the result for the BU setting and Fig. 10a for the BL setting. Interstingly, QPCert seems to 1549 be more tight in a backdoor setting than in a pure poisoning setting. However, this could also be 1550 explained by the fact that an evasion attack is easier to perform than a poisoning attack and thus, 1551 APGD potentially provided lower upper bounds to the actual robustness than for the pure poisoning 1552 setting. Another interesting observation is that for the backdoor settings, the rankings of the GNNs 1553 regarding certified robustness seems to roughly correspond to the robust accuracies obtained by the 1554 backdoor attack.

1555 In Fig. 11 performing a gradient-based attack (APGD) using either exact gradients with QPLayer 1556 or meta-gradients obtained through MetaAttack's surrogate model is compared. For both the PL 1557 (Fig. 11a) and PU (Fig. 11b) setting using exact gradients results in a lower upper bound to the 1558 robust accuracy (i.e., a stronger attack). Thus, we use the exact gradients from QPLayer to measure 1559 the tightness of QPCert. Meta-gradients from MetaAttack are obtained by adapting Algorithm 2 1560 in (Zügner & Günnemann, 2019a) to feature perturbations, through setting a maximum number 1561 of iterations as the stop criterion and instead of choosing an edge with maximal score, update the 1562 feature matrix with the meta-gradient using APGD. In MetaAttack, λ trading of the self-supervised with the training loss is set to 0.5. Interestingly, for small budgets, MetaAttack can lead to the 1563 opposite intended effect. Exemplary, for a GCN in the PL setting with $\delta = 0.1$, the generalization 1564 performance is slightly increased. This indicates that for small perturbation budgets, the meta-gradient 1565 of MetaAttack's surrogate model does not transfer well to the infinite-width networks. However, for

1594

1596

1597

1598

1604

1609

1610

1611

1612

1613

1614

1615

1616

1566 larger budgets, MetaAttack still provides a strong, albeit weaker attack than exact gradients. In Figure 1567 Fig. 12 we compare performing the above mentioned gradient-based backdoor attack with the simple 1568 backdoor strategy proposed by Xing et al. (2024) with a trigger size of 0.5. We observe that Xing 1569 et al. (2024)'s attack is significantly weaker compared to the gradient-based attack and only starts to 1570 reduce accuracy of the models for high attack budgets. This can be explained by several observations: For small ℓ_p -budgets, the backdoor trigger is often distorted in the backdoored nodes by having to 1571 project the perturbation back into the allowed ℓ_p -ball and secondly, the attack is simple, static and not 1572 adaptive. Concretely, it simply copies certain features to other nodes without considering the attacked 1573 model. We want to note that similar to MetaAttack, for small budgets, for BU we can observe for 1574 MLP that the change actually results in slightly higher generalization of the model under attack, 1575 showing that for small budgets, the backdoor strategy in Xing et al. (2024) is not effective. 1576



Figure 9: Tightness of our certificate for data poisoning. Both PU and PL with $p_{adv} = 0.2$ evaluated with APGD (QPLayer).



Figure 10: Tightness of our certificate for backdoor attacks. Both BU and BL again with $p_{adv} = 0.2$.

Robust Accuracy Robust Accuracy 0.8 0.6 0.6 OPCert GIN OPCert GIN MetaAttacl ___. MetaAttack GraphSAGE GCN Skip-PC GraphSAGE GCN GCN Skip-PC GCN APPNP GCN Skip-o APPNP GCN Skip- α 0.4 0.4 MLP MLF SGC SGC 0.1 0.1 0.01 0.02 0.05 0.2 0.50.01 0.02 0.05 0.2 0.5Perturbation Budget $\delta \ (p = \infty)$ Perturbation Budget δ $(p = \infty)$ (a) CSBM: PL (b) CSBM: PU





Figure 12: Comparison of performing a gradient-based backdoor attack versus the simple backdoor attack proposed in Xing et al. (2024).

I.4 Results for p = 2 perturbation budget

1637 We present the results for p = 2 perturbation budget evaluated on CSBM and all the GNNs considered. 1638 We focus on Poison Unlabeled setting. Fig. 13 show the results of the certifiable robustness for all 1639 GNNs and the heatmaps showing the accuracy gain with respect to MLP is in Fig. 13. All the results 1640 are in identical to $p = \infty$ setting and we do not see any discrepancy.



Figure 13: (a): Certifiable robustness for different (G)NNs in Poisoning Unlabeled (PU) for p = 2. (b)-(h): Certified accuracy gain for heatmap for all GNNs. All experiments with Poisoning Unlabeled (PU) and $p_{adv} = 0.2$

1659 1660 1661

1663

1658

1632

1633 1634 1635

1636

1662 I.5 Results for p = 1 perturbation budget

Similar to p = 2, we also present the results for p = 1 perturbation budget evaluated on CSBM and all the GNNs considered for Poison Unlabeled setting in Fig. 14. All the results are in identical to $p = \infty$ setting and we do not see any discrepancy.

- 1668
- 1669
- 1670
- 1671
- 1672
- 1673



Figure 14: (a): Certifiable robustness for different (G)NNs in Poisoning Unlabeled (PU) for p = 1. (b)-(h): Certified accuracy gain for heatmap for all GNNs. All experiments with Poisoning Unlabeled (PU) and $p_{adv} = 0.2$.

1695 I.6 COMPARISON BETWEEN $p = \infty$ and p = 2

1693 1694

1699 1700

1701 1702

1703

1704

1705

1706

1707

1708

1709 1710

1711 1712

1714

1696 1697 We provide a comparison between $p = \infty$ and p = 2 perturbation budget, showing that p = 2 is 1698 tighter than $p = \infty$ for the same budget as expected.



Figure 15: Comparison between $p = \infty$ and p = 2 for Poison Unlabeled setting. $p_{adv} = 0.2$.

1713 I.7 COMPARISON TO COMMON POISONING DEFENSES

In Fig. 16 we compare two common poisoning defenses namely GNNGuard (Zhang & Zitnik, 2020) 1715 and ElasticGNN (Liu et al., 2021) with the certified accuracy provided by QPCert. While the 1716 accuracies provided by a defense are not certified accuracies (i.e., they are only upper bounds to 1717 the true robust accuracy) and hence, can only be compared partly with the certified accuracy which 1718 represents a true lower bound to the robust accuracy. However, a comparison is still interesting as it 1719 allows to answer the question, of how big the gap between the best-certified accuracy to the robust 1720 accuracy provided by defenses is and if we could even get a certified accuracy result comparable to a 1721 poisoning defense's accuracy. Interestingly, Fig. 16 shows that for small to intermediate budgets, the 1722 certified accuracy of an infinite-width GCN as provided by QPCert is higher than the robust accuracy 1723 provided by the defense baselines. This can be explained by the fact that even ElasticGNN and 1724 GNNGuard show lower base clean accuracy despite significant hyperparameter tuning (experimental 1725 details see below paragraph) paired with a few very brittle predictions. We hypothesize that this is due to the difficult learning problem a CSBM poses (despite being a small dataset) paired with the 1726 fact that both poisoning defenses have GCN-like base models where the graph / propagation scheme 1727 is adapted to be more robust to poisoning while potentially trading off clean accuracy.

Both poisoning defenses are trained using the non-negative likelihood loss and the ADAM optimizer following Zhang & Zitnik (2020). GNNGuard uses a 2-layer GCN as a baseline model and hyperpa-rameters are searched in the grid: (i) number of filters $\{8, 16, 32\}$, (ii) dropout $\{0, 0.2, 0.5\}$, (iii) learning rate $\{0.01, 0.001\}$, (iv) weight decay $\{5e-3, 1e-3, 5e-4, 1e-4\}$ over 10 seeds resulting in 720 models. For ElasticGNN the hyperparameter grid reported in Liu et al. (2021) is explored over 10 seeds resulting in 11520 models due to ElasticGNN having more hyperparameters to tune. It's hidden layer size is fixed to 32. Both baseline defenses are attacked using MetaAttack adapted to feature perturbations as done in App. I.3. The infinite-width GCN is attacked using the exact gradient obtained from the QPLayer implementation.



Figure 16: Comparison of different poisoning defenses with the certified accuracy obtained by QPCert.

J ADDITIONAL RESULTS: CORA-MLB

J.1 EVALUATING QPCERT

Fig. 17a shows the certified accuracy on Cora-MLb for the BL settings for $p_{cert} = 0.1$. Figs. 17b to 17d and 18 show a detailed analysis into the certified accuracy difference of different GNN architectures for PU setting for $p_{cert} = 0.1$.



1768Figure 17: (a) Backdoor Labeled (BL) Setting. (b)-(d) Heatmaps of GCN, SGC, and APPNP for1769Poison Unlabeled (PU) setting on Cora-MLb with $p_{adv} = 0.1$.1770









Fig. 19 shows that the inflection point observed in Fig. 4c is not observed in the other settings.

Figure 20: Influence of symmetric and row normalized adjacency in GCN and SGC for poison unlabeled and poison labeled settings. 1812

Results on p = 1 adversary J.4

1813 1814

1815

1826

1827 1828 1829

1830

1816 Fig. 21 shows the certifiable robustness to p = 1 adversary on Cora-MLb dataset. The observation is 1817 consistent to the CSBM case.



Figure 21: Cora-MLb results for PL, PU, BL and BU under p = 1 perturbation.

Κ ADDITIONAL RESULTS: CORA-ML

1831 For Cora-ML we choose 100 test nodes at random and investigate in Fig. 22a the poison labeled (PL) setting with a strong adversary $p_{adv} = 1.0$ for GCN, SGC and MLP. It shows that QPCert can 1833 provide non-trivial robustness guarantees even in multiclass settings. Fig. 22b shows the results for poison unlabeled (PU) and $p_{adv} = 0.05$. Only SGC shows better worst-case robustness than MLP. 1834 This, together with both plots showing that the certified radii are lower compared to the binary-case, 1835 highlights that white-box certification of (G)NNs for the multiclass case is a more challenging task.







Figure 25: Heatmaps of GCN Skip- α , GCN Skippc, GraphSAGE, and GIN for Poison Unlabeled (*PU*) setting on WikiCSb with $p_{adv} = 0.02$.

1904 M RELATED WORKS

1901

1902 1903

1905 **Poisoning certificates.** The literature on poisoning certificates is significantly less developed than certifying against test-time (evasion) attacks (Li et al., 2023) and we provide an overview in Table 1. 1907 Black-box certificates for poisoning are derived following three different approaches: (i) Randomized 1908 smoothing, a popular probabilistic test-time certificate strategy (Cohen et al., 2019), in which 1909 randomization performed over the training dataset (Rosenfeld et al., 2020; Weber et al., 2023; Zhang 1910 et al., 2022). Other than data partitioning, a common defense is to sanitize the data, and Hong 1911 et al. (2024) certifies diffusion-based data sanitation via randomized smoothing. (*ii*) Ensembles: 1912 Creating separate partitions of the training data, training individual base classifiers on top of them and 1913 certifying a constructed ensemble classifier (Levine & Feizi, 2021; Jia et al., 2021; Wang et al., 2022; Rezaei et al., 2023); Jia et al. (2021) and Chen et al. (2022) offer certificates and collective certificates, 1914 respectively, for bagging, while Levine & Feizi (2021) and Wang et al. (2022) derive certificates for 1915 aggregation-based methods tailored for black-box classifiers. (*iii*) Differential Privacy² (DP): Ma 1916 et al. (2019) show that any (ϵ, δ) -DP learner enjoys a certain provable poisoning robustness. Liu et al. 1917 (2023) extend this result to more general notions of DP. Xie et al. (2023) derives guarantees against 1918 arbitrary data poisoning in DP federated learning setup. However, white-box deterministic poisoning 1919 certificates remain sparse. Drews et al. (2020) and Meyer et al. (2021) derive poisoning certificates 1920 for decision trees using abstract interpretations, while Jia et al. (2022) provides a poisoning certificate 1921 for nearest neighbor algorithms based on their inherent majority voting principle. Recently, Bian et al. 1922 (2024) derives a poisoning certificate for naive Bayes classification. 1923

Poisoning attacks and defense using the bilevel problem. Biggio et al. (2012) and Xiao et al. (2015)
use the bilevel problem with SVM hinge loss and regularized ERM to generate poison samples, and solve it using iterative gradient ascent. Mei & Zhu (2015) also recognize the possibility to transform the bilevel problem into a single-level one. However, they only reformulate the problem into a single-level bilinear one and solve it heuristically w.r.t. to the training data to generate poisoning attacks. Koh & Liang (2017) also considers the bilevel problem to detect and also generate poisoned samples using influence functions (gradient and Hessian vector product).

Graphs. Currently, there are no poisoning certificates for clean-label attacks specifically developed for GNNs or the task of node classification. (Lai et al., 2024) is the only work on poisoning certification of GNNs, but differ incomparably in their threat model and are black-box as well as not applicable to backdoors. However, there are many works on certifying against test-time attacks on graphs and Günnemann (2022) provides an overview.

- 1935
- 1330
- 1938
- 1939
- 1940
- 1941
- 1942 1943

²The mechanism to derive a poisoning certificate from a certain privacy guarantee is model agnostic, thus we count it as black-box. However, the calculated privacy guarantees may depend on white-box knowledge.

¹⁹⁴⁴ N FURTHER DISCUSSIONS

1946 N.1 Applicability to Commonly Studied Perturbation Models and Attacks

1948 QPCert applies to any poisoning or backdoor attack that performs ℓ_p -bounded feature perturbations. As such, QPCert is directly applicable to clean-label (graph) backdoor attacks as proposed by Turner 1949 et al. (2019) and Xing et al. (2024), and clean-label poisoning attacks such as Huang et al. (2020) and 1950 Geiping et al. (2021). It is not directly applicable to poisoning of the graph structure as performed 1951 by MetaAttack (Zügner & Günnemann, 2019a). However, the MetaAttack strategy can be easily 1952 adapted to poison node features as done in App. I.3 and we discuss the challenges to extend QPCert 1953 to structure perturbations in App. N.2. The backdoor attack proposed by Dai et al. (2023) changes 1954 the node features and training labels jointly and thus, our method is only applicable if the training 1955 labels will be kept constant or the poisoned nodes are sampled only from the class, the training label 1956 should be changed to. Similar, Xi et al. (2021) develops a backdoor attack that changes the features 1957 and graph structure jointly and thus, QPCert is not applicable given the graph structure changes. 1958

1959 N.2 CERTIFYING AGAINST GRAPH STRUCTURE PERTURBATIONS 1960

1961 In the following, we discuss how to approach certifying against poisoning of the graph structure and 1962 the open challenges that arise in the process. To certify against poisoning the graph structure, again 1963 Eq. (3) has to be solved but now, the adversary \mathcal{A} can change the graph structure instead of the node 1964 features, meaning the optimization in Eq. (3) is performed w.r.t. the graph structure matrix S and 1965 thus, reads

1966 1967

$$\min_{\widetilde{\boldsymbol{S}},\theta} \mathcal{L}_{att}(\theta,\widetilde{\mathcal{G}}) \quad \text{s. t.} \quad \widetilde{\boldsymbol{S}} \in \mathcal{A}(\boldsymbol{S}) \land \theta \in \operatorname*{arg\,min}_{\theta'} \mathcal{L}(\theta',\widetilde{\mathcal{G}})$$
(37)

1968 1969

with $\widetilde{S} \in \mathcal{A}(S)$ representing the perturbed graph structure matrices constructable by the adversary and $\widetilde{\mathcal{G}} = (\widetilde{S}, \mathbf{X})$. Indeed, it will be possible to reformulate this problem into a single-level problem similar to the description in Sec. 3. While in theory, QPCert Theorem 1 also applies to structure perturbations, the bounding strategy from Sec. 3.1 does result in loose bounds for structure perturbations, as an untrusted node will always result in a lower bound in the respective adjacency matrix entry of 0 and an upper bound of 1 - thus, spanning the whole space of possible entries.

To overcome this, one can approach certifying graph structure perturbations by including the NTK computation into the optimization problem with the drawback that each type of GNN architecture will require slight adaptations of the optimization problem depending on its corresponding NTK. Assuming that the chosen model is an L = 1 layer GCN (the formulation can be easily extended to arbitrary layers, see App. B.3), the bilevel optimization problem reads as follows

1989 1990

$$\min_{\boldsymbol{\alpha}, \widetilde{\boldsymbol{S}}, \boldsymbol{Q}, \boldsymbol{\Sigma}_1, \boldsymbol{\Sigma}_2, \boldsymbol{E}_1, \dot{\boldsymbol{E}}_1, \dot{\boldsymbol{E}}_2} \operatorname{sgn}(\hat{p}_t) \sum_{i=1}^m y_i \alpha_i Q_{ti} \quad s.t. \quad \widetilde{\boldsymbol{S}} \in \mathcal{A}(\boldsymbol{S}) \land \boldsymbol{\alpha} \in \mathcal{S}(\boldsymbol{Q})$$
(38)

$$\boldsymbol{Q} = \widetilde{\boldsymbol{S}}(\boldsymbol{\Sigma}_1 \odot \dot{\boldsymbol{E}}_1) \widetilde{\boldsymbol{S}}^T + \boldsymbol{\Sigma}_2 \odot \dot{\boldsymbol{E}}_2 \qquad (39)$$

$$\boldsymbol{\Sigma}_1 = \widetilde{\boldsymbol{S}} \boldsymbol{X} \boldsymbol{X} \widetilde{\boldsymbol{S}}^T \tag{40}$$

$$\boldsymbol{\Sigma}_2 = \boldsymbol{S} \boldsymbol{E}_1 \boldsymbol{S}^T \tag{41}$$

$$\mathbf{E}_{1} = c_{\sigma} \underset{\mathbf{F} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_{l})}{\mathbb{E}} \left[\sigma(\mathbf{F}) \sigma(\mathbf{F})^{T} \right]$$
(42)

1991
1991

$$\mathbf{\dot{E}}_{1} = c_{\sigma} \underset{\mathbf{F} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_{l})}{\mathbb{E}} \begin{bmatrix} \dot{\sigma}(\mathbf{F}) \dot{\sigma}(\mathbf{F})^{T} \end{bmatrix}$$
 (43)
1993

$$\dot{\mathbf{E}}_2 = \mathbf{1}_{n \times n} \tag{44}$$

1996 This (non-linear) bilevel problem can be transformed into a single-level problem as described in 1997 Sec. 3, as the inner-level problem $\alpha \in S(Q)$ is the same as in Eq. (4) and the same strategy can be applied to linearly model the resulting constraints from the KKT conditions. However, a crucial

2004

2005 2006 2007

20

2015

2022

2023

2034

2035

1998 difference to Eq. (4) are the additional non-linear constraints arising from optimizing over the NTK computation. Eqs. (39) to (41) are multilinear constraints that can be reduced to bilinear constraints 2000 by introducing additional variables as follows (for brevity, again writing the problem in its bilevel 2001 form):

> $\min \ \mathrm{sgn}(\hat{p}_t) \sum_{i=1}^m y_i \alpha_i Q_{ti} \quad s.t. \quad \widetilde{\boldsymbol{S}} \in \mathcal{A}(\boldsymbol{S}) \ \land \ \boldsymbol{\alpha} \in \mathcal{S}(\boldsymbol{Q})$ (45)

$$\boldsymbol{Q} = \boldsymbol{H}_1^{\prime\prime} + \boldsymbol{H}_2 \tag{46}$$

$$H_1 = \Sigma_1 \odot \dot{E}_1 \tag{47}$$

2009

$$H_1' = H_1 S^T$$
 (48)

 2010
 $H_1'' = SH_1'$ (49)

 2011
 $H_1'' = SH_1'$ (49)

$$H_2 = \Sigma_2 \odot \dot{E}_2 \tag{50}$$

2013
$$\Sigma_1 = M_1 M_1^T$$
(51)
2014
$$M_2 = S X$$
(52)

$$M_2 = E_1 S^T$$
(53)

2017
$$\Sigma_2 = SM_2$$
(54)
2018
$$\mathbf{E}_1 = c_{\sigma} \mathop{\mathbb{E}}_{\mathbf{F} \sim \mathcal{N}(\mathbf{0}, \mathbf{\Sigma}_1)} \left[\sigma(\mathbf{F}) \sigma(\mathbf{F})^T \right]$$
(55)

2020

$$\dot{\mathbf{E}}_{1} = c_{\sigma} \underset{\mathbf{F} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_{1})}{\mathbb{E}} \left[\dot{\sigma}(\mathbf{F}) \dot{\sigma}(\mathbf{F})^{T} \right]$$
(56)

$$\dot{\mathbf{E}}_2 = \mathbf{1}_{n \times n} \tag{57}$$

where the optimization is over the same variables as in the previous problem, and additionally the 2024 variables $H_1, H_1', H_1'', H_2, M_1$, and M_2 . Eqs. (46) and (52) are linear constraints, the rest of 2025 the newly introduced constraints represent bilinear terms. The same bilinearization strategy can be 2026 applied given the NTK computation over arbitrary layers. The remaining non-linear and non-bilinear 2027 terms are Eqs. (55) and (56). They can be solved in closed-form resulting in relatively well-behaved 2028 functions as shown in App. D.1. Thus, a convex relaxation of the expectation terms can be derived by 2029 e.g., choosing linear functions that lower and upper bound the expectation terms. 2030

Assume for now that one can linearly model $S \in \mathcal{A}(S)$, this can be achieved by e.g., choosing the 2031 adjacency matrix without normalization as graph structure matrix as done by Hojny et al. (2024). 2032 Then, the crucial question is: 2033

How to effectively solve the arising bilinear problem?

2036 In particular, the bilinearities arise in both, the constraints and objective, and thereby this contrasts 2037 e.g., with Zügner & Günnemann (2020) who only have to deal with a bilinear objective but have 2038 linear constraints. The problem can be slightly simplified if S is chosen to be the unnormalized 2039 adjacency matrix, as then any S is discrete and thus Eqs. (48), (49) and (53) represent multiplications 2040 of a continuous with a discrete variable and thus, can be linearly modeled using standard modeling 2041 techniques. However, the objective and Eqs. (47), (50) and (51) remain products of continuous variables and thus, can fundamentally not be modeled linearly. One potential way to tackle this, is to 2042 use techniques of convex relaxations of bilinear functions as e.g., the so called McCormic envelope 2043 (McCormick, 1976). However, it is not clear if common bilinear relaxation techniques can scale to 2044 problems of the size necessary to compute practical certificates for machine learning datasets, nor 2045 is it clear if the relaxations introduced in the process result in tight enough formulations to yield 2046 non-trivial certificates. This is complicated by the fact that problems that are studied in the bilinear 2047 optimization literature are often significantly smaller than the problem size we can expect from 2048 certifying graph structure perturbations. However, it is not unlikely that further progress in bilinear 2049 optimization will make this problem tractable. 2050

We want to note that linearly modeling $S \in \mathcal{A}(S)$ by choosing an unnormalized adjacency matrix as 2051 the graph structure matrix results in a restriction of possible architecture to certify. It could be possible to adapt Zügner & Günnemann (2020)'s modeling technique for the symmetric-degree normalized adjacency matrix they use to certify a finite-width GCN to the above optimization problem without increasing its difficulty as it is already bilinear.

2055 2056

2057 N.3 PRACTICAL IMPLICATION OF FEATURE AND STRUCTURE PERTURBATIONS

2058 Both feature and structure perturbations can find applications in real-world scenarios. In particular, 2059 important application areas for graph learning methods with adversarial actors are fake news detection 2060 (Hu et al., 2024) and spam detection (Li et al., 2019). Regarding fake news detection, feature 2061 perturbations can model changes to the fake news content or to (controlled) user account comments and profiles to mislead detectors (Hu et al., 2024; Le et al., 2020). Structure perturbations allow to 2062 model a change in the propagation patterns (e.g., through changing a retweet graph) (Wang et al., 2063 2023). The qualitative difference in the application of feature compared to structure perturbations 2064 is similar for spam detection. Here, feature perturbation can model spammers trying to adapt their 2065 comments to avoid detection (Li et al., 2019; Wang et al., 2012). However, structure perturbations 2066 can model behavioral changes in the posting patterns of spammers to imitate real users (Soliman & 2067 Girdzijauskas, 2017; Wang et al., 2012). 2068

2069 2070

2076

2077

2078

2079

2080

2081

2082

2083

2084

2085 2086

2087

N.4 QPCERT FOR OTHER GNNS

While our analysis focused on commonly used GNNs with and without skip connections, QPCert is broadly applicable to any GNNs with a well-defined analytical form of the NTK. The following challenges and considerations have to be taken into account when extending QPCert to other architectures:

- 1. **NTK-network equivalence:** The equivalence between the network and NTK breaks down when the network has non-linear last layer or bottleneck layers (Liu et al., 2020). Consequently, our certificates do not hold for such networks.
- 2. **Analytical form of NTK:**Deriving a closed-form expression for NTK is needed to derive bounds on the kernel. This might be challenging for networks with batch-normalizations or advanced pooling layers.
 - 3. **Tight bounds for the NTK:** Ensuring non-trivial certificates requires deriving tight bounds for the NTK. Depending on the NTK, additional adaptation of our bounding strategy may be necessary.

Despite these considerations, most message-passing networks satisfy these criteria, making QPCert readily applicable to a wide range of architectures.

2089 2090 N.5 Adaptation to QPCert for Graph Classification

Our work can be extended to graph classification using the graph NTK of the corresponding neural network trained for the task. Note that in this case, the kernel is computed between all pairwise graphs instead of nodes. Du et al. (2019) derived one such NTK for graph classification. Using this NTK and our MILP, robustness certificate can be derived. We elaborate on the technical details below.

Adversary. First, we extend our adversary setting to include multiple graphs and allow for node feature perturbation. We have n graphs $\mathcal{G}_i = (\mathbf{S}_i, \mathbf{X}_i), \mathbf{S}_i \in \mathbb{R}^{n_i \times n_i}, \mathbf{X}_i \in \mathbb{R}^{n_i \times d}$ where n_i is the number of nodes in graph i with d dimensional features, for all $i \in [n]$ and the adversary \mathcal{A} can perturb the features \mathbf{X} where $\in \mathcal{B}_p(\mathbf{x})$. As we consider semi-verified learning setting, let \mathcal{U}_i be the set of nodes in \mathcal{G}_i that are potentially controlled by the adversary \mathcal{A} .

2100 2101 MILP. The certification framework applies directly without any change given the NTK Q and its 2102 element-wise bounds. Therefore, the important adaptation here is to derive bounds on the NTK. 2103 Since, in this setting, the NTK computation involves all pairwise feature matrix covariance, that is, 2104 $\mathbf{X}_i \mathbf{X}_j^T$ for all i, j pairs. Similar to node classification setting, we consider $\mathbf{X}_i \in \mathcal{A}(\mathbf{X}_i)$ and derive 2105 $\mathbf{X}_i \mathbf{X}_j^T = \mathbf{X}_i \mathbf{X}_j^T + \Delta_{ij}$. We derive the bounds for Δ_{ij} for the perturbation $\mathcal{B}_p(\mathbf{x})$ with p = 2. It is 2106 easy to extend to p = 1 and $p = \infty$ in a similar way. Bounds for Δ_{ij} and p = 2. We consider the perturbed feature matrix $\widetilde{\mathbf{X}}_i \in \mathcal{A}(\mathbf{X}_i)$ for all graphs *i*, and $\widetilde{\mathbf{X}}_i = \mathbf{X}_i + \mathbf{\Gamma}_i \in \mathbb{R}^{n_i \times d}$ where $\mathbf{\Gamma}_{ij}$ is the adversarial perturbations added to node *j* of graph *i* by the adversary, therefore, $\|\mathbf{\Gamma}_{ij}\|_p \leq \delta$ and $\mathbf{\Gamma}_{ij} > 0$ for $j \in \mathcal{U}_i$ and $\mathbf{\Gamma}_{ij} = 0$ for $j \notin \mathcal{U}_i$. Then

2110
$$\widetilde{\mathbf{X}}_{i}\widetilde{\mathbf{X}}_{i}^{T} = (\mathbf{X}_{i} + \boldsymbol{\Gamma}_{i})(\mathbf{X}_{i} + \boldsymbol{\Gamma}_{i})^{T}$$

$$= \mathbf{X}_i \mathbf{X}_j^T + \mathbf{\Gamma}_i \mathbf{X}_j^T + \mathbf{X}_i \mathbf{\Gamma}_j^T + \mathbf{\Gamma}_i \mathbf{\Gamma}_j^T = \mathbf{X}_i \mathbf{X}_j^T + \Delta_{ij}.$$
(58)

2114 As a result, it suffices to derive the element-wise worst-case bounds for $\Delta_{ij} \in \mathbb{R}^{n_i \times n_j}$, $(\Delta_{ij})_{ab}^L \leq$ 2115 $(\Delta_{ij})_{ab} \leq (\Delta_{ij})_{ab}^U$, for different perturbations. To do so, our strategy is to bound the scalar products 2116 $\langle (\Gamma_i)_a, (\mathbf{X}_j)_b \rangle$ and $\langle (\Gamma_i)_a, (\Gamma_j)_b \rangle$ element-wise.

$$\begin{aligned} |\langle (\mathbf{\Gamma}_i)_a, (\mathbf{X}_j)_b \rangle| &\leq \|(\mathbf{\Gamma}_i)_a\|_2 \|(\mathbf{X}_j)_b\|_2 \leq \delta \|(\mathbf{X}_j)_b\|_2\\ |\langle (\mathbf{\Gamma}_i)_a, (\mathbf{\Gamma}_j)_b \rangle| &\leq \|(\mathbf{\Gamma}_i)_a\|_2 (\|\mathbf{\Gamma}_j)_b\|_2 \leq \delta^2. \end{aligned}$$
(59)

2121 Using Eq. (59), we derive the lower and upper bounds of $(\Delta_{ij})_{ab}$:

$$\begin{array}{l} \textbf{2122} \\ \textbf{2123} \\ \textbf{2124} \\ \textbf{2124} \\ \textbf{2125} \\ \textbf{2126} \end{array} \quad (\Delta_{ij})_{ab}^{L} = \begin{cases} 0+ & \text{if } a \notin \mathcal{U}_{\flat}, b \notin \mathcal{U}_{j} \\ -\delta ||(\mathbf{X}_{j})_{b}||_{2} + & \text{if } a \in \mathcal{U}_{i} \\ -\delta ||(\mathbf{X}_{i})_{a}||_{2} + & \text{if } b \in \mathcal{U}_{j} \\ -\delta^{2} & \text{if } a \in \mathcal{U}_{i}, b \in \mathcal{U}_{j} \end{cases} \quad (\Delta_{ij})_{ab}^{L} = \begin{cases} 0+ & \text{if } a \notin \mathcal{U}_{\flat}, b \notin \mathcal{U}_{j} \\ \delta ||(\mathbf{X}_{j})_{b}||_{2} + & \text{if } a \in \mathcal{U}_{i} \\ \delta ||(\mathbf{X}_{i})_{a}||_{2} + & \text{if } b \in \mathcal{U}_{j} \\ \delta^{2} & \text{if } a \in \mathcal{U}_{i}, b \in \mathcal{U}_{j} \end{cases}$$

2128 Thus, writing the bounds succinctly using the indicator function gives us,

$$(\Delta_{ij})_{ab}^{L} = -\delta \|(\mathbf{X}_{j})_{b}\|_{2} \mathbb{1}[a \in \mathcal{U}_{i}] - \delta \|(\mathbf{X}_{i})_{a}\|_{2} \mathbb{1}[b \in \mathcal{U}_{j}] - \delta^{2} \mathbb{1}[a \in \mathcal{U}_{i} \land b \in \mathcal{U}_{j} \land a \neq b],$$

$$(\Delta_{ij})_{ab}^{U} = \delta \|(\mathbf{X}_{j})_{b}\|_{2} \mathbb{1}[a \in \mathcal{U}_{i}] + \delta \|(\mathbf{X}_{i})_{a}\|_{2} \mathbb{1}[b \in \mathcal{U}_{j}] + \delta^{2} \mathbb{1}[a \in \mathcal{U}_{i} \land b \in \mathcal{U}_{j} \land a \neq b].$$

2132
2133 Using these bounds, similar to node classification, we can propagate them through the NTK computa2134 tion to get the bounds on NTK. Using the NTK bounds, we can apply QPCert to get the certificate for
2135 graph classification. □