

EVA: EVOLUTIONARY ATTACKS ON GRAPHS

Anonymous authors

Paper under double-blind review

ABSTRACT

Even a slight perturbation in the graph structure can cause a significant drop in the accuracy of graph neural networks (GNNs). Most existing attacks leverage gradient information to perturb edges. This relaxes the attack’s optimization problem from a discrete to a continuous space, resulting in solutions far from optimal. It also restricts the adaptability of the attack to non-differentiable objectives. Instead, we propose an evolutionary-based algorithm to solve the discrete optimization problem directly. Our Evolutionary Attack (EvA) works with any black-box model and objective, eliminating the need for a differentiable proxy loss. This permits us to design two novel attacks that: reduce the effectiveness of robustness certificates and break conformal sets. We introduce a sparse encoding that results in memory complexity that is linear in the attack budget. EvA reduces the accuracy by an additional $\sim 11\%$ on average compared to the best previous attack, revealing significant untapped potential in designing attacks.

1 INTRODUCTION

Given the widespread applications of graph neural networks (GNNs), studying their robustness to natural and adversarial noise is of great importance. In node classification, GNNs leverage the edge structure between data points to improve their performance. However, a small perturbation in the graph structure (adding or removing a few edges) can significantly reduce GNNs’ accuracy, even below the performance of an MLP (which completely discards the structure). Similar to images and continuous data, most of the proposed (structure) attacks are gradient-based. They compute the gradients of a loss w.r.t. the adjacency matrix and apply a perturbation according to that. Gradient-based attacks face several challenges. They solve a relaxation of the original combinatorial (discrete) optimization problem – the entries of the adjacency matrix are relaxed from $\{0, 1\}$ to $[0, 1]$. They need a differentiable proxy loss function since the actual objective of the attacker (e.g. accuracy) is often not differentiable, and the usual proxy such cross-entropy is suboptimal (Geisler et al., 2023). They assume white-box access to the model, including the structure and the weights. This limits the applicability or requires surrogate models. They can provide a false sense of security since defenses may be obfuscating gradients (Athalye et al., 2018; Geisler et al., 2023) and can get stuck in local minima. Their memory complexity grows quadratically w.r.t the number of nodes. Although the adjacency matrix is often sparse, the gradients w.r.t. it are not. As a result, tricks like block coordinate descent are needed (Geisler et al., 2021). We propose a model-agnostic evolutionary attack (EvA) that fixes all five of the above issues.

EvA explores the space of possible perturbations with a genetic algorithm (GA). Our approach operates in the discrete space of potential perturbations without information from gradients – avoiding relaxation. It directly optimizes the objective (like accuracy) as long as it provides a meaningful signal. In addition to eliminating the need for a differentiable proxy, this black-box access to the objective enables us to define a broader class of attacks. In fact, it allowed us to easily design two novel attacks on graphs that aim at decreasing the effectiveness of robustness certificates or that break conformal guarantees. EvA shows outstanding effectiveness on vanilla and adversarially trained models compared to SOTA attacks. Unlike the gradient-based attacks, our attack has a $\mathcal{O}(\epsilon \cdot E)$ memory complexity where ϵ is the perturbation budget, and E is the number of edges. This is because instead of storing a squared block of gradients, which scales with the size of the adjacency matrix, we only store the edge perturbations as an index. During the evaluation, we also maintain the same sparsity in representation as the graph itself. To take advantage of the available free memory, we employ a batch evaluation approach that speeds up the optimization.

Adversarial attacks are supposed to be imperceptible. In images, this is modeled by a L_p -ball of a small radius. Similarly, for the graph structure, a commonly used metric is the L_0 ball, which allows changing of the node degree significantly. Since this may be perceptible, we can incorporate *constraints* in our attack that limit the number of perturbations per node, in addition to the global budget. Similar to gradient-based attacks (Geisler et al., 2021), we set this so-called local budget to a fraction of the node’s original degree. Interestingly, in some cases, our constrained attack can even beat the best unconstrained gradient-based attack. Overall, EvA finds significantly better solutions compared to the previous state-of-the-art methods (Geisler et al., 2021; Gosch et al., 2024), which highlight the sub-optimality of gradient-based methods.

Given the black-box nature of EvA, we were easily able to introduce the first graph certificate attack. One defense against adversarial attacks is to certify the prediction of a (smoothed) classifier (Bojchevski et al., 2020).

Certificates provide a robustness guarantee that the prediction will not change given a limited set of possible perturbations (e.g., at most r_a additions and r_d deletions). The certified ratio is the fraction of nodes for which the guarantee holds. Here, we define the attacker’s objective as decreasing the certified ratio. EvA can decrease the ratio below the MLP level (which is by definition robust to any perturbation in structure), while also preserving the clean accuracy – making it less noticeable to a defender. We also introduce the first conformal attack on graphs. Conformal prediction (CP) converts any model’s output to prediction sets with a guarantee to cover the true label with (adjustable) high probability. With EvA we can attack these conformal sets to either break the guarantee of increase the sets size (making them useless). While in principle one can design gradient-based attacks for these two new objectives, the amount of work is nontrivial since there are many non-differentiable components that would need to be relaxed. In contrast, for EvA, designing a new attack is simply a matter of changing the fitness function of the GA.

Importantly, perhaps the main contributions of this work is to highlight a scarcely explored research direction for attacks. Even off-the-shelf genetic algorithms significantly outperform gradient-based attacks. Fig. 1 compares EvA to the other attacks proposed over time. Our custom adaptive mutation further improves performance, but we argue that the space of evolutionary (and more broadly search-based) attacks has a lot of untapped potential.

2 BACKGROUND AND RELATED WORK

Problem setup. We focus on attacking the semi-supervised node classification task on graphs via perturbing a small number of edges. Formally, we are given a graph $\mathcal{G} = (\mathbf{X}, \mathbf{A}, \mathbf{y})$ in which \mathbf{X} is the features matrix assigning a feature vector x_i to each node v_i in the graph, \mathbf{A} is the adjacency matrix (often sparse) that represents the set of edges \mathcal{E} , and \mathbf{y} is the partially observable vector of labels. Nodes are partitioned into labeled and unlabeled sets $\mathcal{V} = \mathcal{V}_l \cup \mathcal{V}_u$. The GNN is trained on an observed subgraph \mathcal{G}_{tr} that includes the labeled nodes. Gosch et al. (2024) argue that the transductive setup, where $\mathcal{G}_{tr} = \mathcal{G}$ is unrealistic since perfect robustness can be achieved by memorizing the training graph. Therefore, we mainly focus on the inductive setting where a model f is trained on an induced subgraph $\mathcal{G}_{tr} \subseteq \mathcal{G}$, validated on $\mathcal{G}_{val} \subseteq \mathcal{G}$ and tested on \mathcal{G}_{test} where $\mathcal{G}_{tr} \subset \mathcal{G}_{val} \subset \mathcal{G}_{test} = \mathcal{G}$.

Threat model. Our goal is to find a perturbation matrix $\mathbf{P} \in \{0, 1\}^{n \times n}$ that flips entities of the adjacency matrix $\hat{\mathbf{A}} = \mathbf{A} \oplus \mathbf{P}$ to decrease the accuracy as much as possible. Here $n = |\mathcal{V}|$, and \oplus is the element-wise XOR operator. For a given function f as the GNN model, the accuracy is defined as $\sum_{v_i \in \mathcal{V}_{att}} (1/|\mathcal{V}_{att}|) \cdot \mathbf{1}[f(\mathcal{G})_{v_i} = y_i]$ where \mathcal{V}_{att} is the set of nodes that we attack. In global attacks this is usually the test nodes, while in targeted attacks the target is a single node. To keep the perturbations imperceptible, we assume that the adversary can only perturb up to $\delta := \epsilon \cdot |\mathcal{E}[\mathcal{V}_{att} : \mathcal{V}]|$ edges where $\mathcal{E}[\mathcal{A} : \mathcal{B}]$ is the subset of edges between nodes in \mathcal{A} and \mathcal{B} . Formally, for any generic

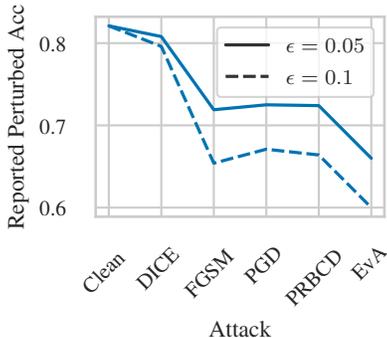


Figure 1: Reported performance of various attacks for transductive setting

108 loss function \mathcal{L} ,

$$109 \quad \mathbf{P} = \arg \max_{\mathbf{P}} \mathcal{L}(f(\mathcal{G}(\mathbf{X}, \mathbf{A} \oplus \mathbf{P}))_{\text{att}}, \mathbf{y}_{\text{att}}) \\ 110 \quad \text{s.t.} \quad \mathbf{1}_N \mathbf{P} \mathbf{1}_N^\top \leq \epsilon \cdot |\mathcal{E}[\mathcal{V}_{\text{att}} : \mathcal{V}]| \quad (1)$$

111 Here $f(\cdot)_{\text{att}}$ returns the vector of predictions for the nodes in \mathcal{V}_{att} . In an evasion attack, \mathcal{L} is the
112 accuracy. Eq. 1 can include additional constraints like the local constraint from Gosch et al. (2023)
113 that restrict the number of perturbation per node to some fraction (e.g., half) of its degree.

114 **Gradient-based attacks.** A common approach to attack the graph structure is to compute the gradient
115 of the loss function w.r.t. the adjacency matrix. This requires a relaxation on the domain of \mathbf{A} from
116 $\{0, 1\}^{n \times n}$ to $[0, 1]^{n \times n}$. If the loss function is not differentiable (e.g., accuracy), then a differentiable
117 surrogate like the categorical cross entropy or tanh-margin (Geisler et al., 2023) is used instead.
118 We compute the derivatives of the loss w.r.t. \mathbf{A} and update the perturbation matrix. Finally, the
119 edges are either sampled or rounded from the perturbation matrix, which returns the solution to the
120 binary domain. There are various tricks to improve gradient-based attacks, but most follow a similar
121 high-level procedure.

122 **Related Work.** Adversarial attacks on graphs are generally divided into two main categories: evasion
123 attacks Xu et al. (2019); Zügner et al. (2018); Geisler et al. (2023); Gosch et al. (2024), where the
124 attacker perturbs the graph after the model has been trained, and poisoning attacks Zügner et al.
125 (2020); Lingam et al. (2023); Zügner et al. (2018), where the attacker modifies the graph prior to
126 training. These attacks can be further classified into global attacks (e.g., Geisler et al. (2023); Zhu
127 et al. (2023)), which target multiple node predictions simultaneously, and targeted attacks, which
128 focus on a single node or a subset of nodes. The manipulations can involve altering node attributes,
129 modifying edge structures, or introducing malicious nodes. The earliest adversarial attacks on graphs
130 were inspired by techniques used on continuous data, utilizing gradients to approximate perturbations
131 on inherently discrete edges Xu et al. (2019); Zügner et al. (2018); Geisler et al. (2023). Additionally,
132 reinforcement learning has been employed as an alternative approach to execute adversarial attacks
133 Dai et al. (2018). Attackers leverage reinforcement learning algorithms to refine their attack strategies
134 and disrupt the learning process of GNNs Sun et al. (2023). Although some new attacks have been
135 proposed in recent years (e.g., by Zhang et al. (2024; 2023); Wang et al. (2023)), they are all based
136 on some traditional algorithms (like gradient-based methods).

137 3 EVA: EVOLUTIONARY ATTACK

138 Our evolutionary-based attack (EvA) uses a genetic algorithm (Holland, 1984) as a heuristic to directly
139 optimize Eq. 1. We define an initial set of possible (candidate) perturbations – called “population” –
140 and iteratively improve this population. In each iteration, the individuals in the population are ordered
141 based on their fitness, specifically in terms of how much each individual decreases the accuracy. We
142 draft the next population by keeping the best individuals and producing new ones as a function of
143 them. Our population for the next iteration is finalized after a mutation which introduces additional
144 randomness that helps with exploration. Each element in the population is a possible perturbation,
145 which is encoded as a vector of indices where an edge is flipped.

146 **Genetic algorithm (GA) in EvA.** We can define a genetic solver through the definition of four
147 main components. (i) Population: It is a set of feasible answers to the problem which gradually
148 improve over iterations. In our case each population element is one potential perturbation on the
149 adjacency matrix. We define mapping $\Pi : \mathbf{x}_{i,t} \in [\frac{n}{2}(n-1)]^\delta \mapsto [n]^2$ which is an enumeration
150 on the upper triangle of the $n \times n$ adjacency matrix. With that, we define each candidate as set
151 $\mathbf{s}_{i,t} \in [\frac{n}{2}(n-1)]^\delta$ which refers to a perturbation. The corresponding perturbation matrix $\mathbf{P}_{i,t}$ is
152 simply defined as $\mathbf{P}_{i,t}[p, q] = \mathbf{P}_{i,t}[q, p] = 1 \Leftrightarrow \exists j : \mathbf{s}_{i,t}[j] = \Pi^{-1}(p, q)$ for $p < q$. (ii) Fitness:
153 Is a notion of how close to optimal each population element is. Given any loss function \mathcal{L} we
154 define the fitness function $\text{fit} : [\frac{n}{2}(n-1)]^\delta \mapsto \mathbb{R}$, as $\text{fit}(\mathbf{s}) = \mathcal{L}(\mathbf{X}, \mathbf{A} \oplus \mathbf{P}_{\mathbf{s}}, \mathbf{y})$. Note that this
155 objective can be non-differentiable, such as accuracy. As long as the loss function has enough
156 sensitivity to differentiate between various individuals, we use it directly as the fitness (see § 4 for
157 extended discussion). (iii) Crossover: Is an operation that defines a new population element by
158 combining two existing ones. The crossover operation at point j defines a new candidate vector
159 $\mathbf{s}_{\text{new}} = \text{cross}_j(\mathbf{s}_1, \mathbf{s}_2) := \mathbf{s}_1[:j] \bullet \mathbf{s}_2[j+1:]$ where \bullet is the concatenation of two vectors. Crossover
160

operation with more than one point is defined recursively in the order of joints. The number of crossovers k_{cross} is a hyperparameter (see § C), and their location is chosen randomly in the range of the perturbation size. (iv) Mutation: **Is a random operation that allows further exploration.** The function $\text{mutate} : [\frac{n}{2}(n-1)]^\delta \mapsto [\frac{n}{2}(n-1)]^\delta$ is a random mapping of a candidate to another. One simple mutation function changes each index with some mutation probability p to some other index in the range (uniformly at random). In § 4 we discuss more advanced mutation strategies that significantly improve performance.

Given all the ingredients above, GA operates by iteratively evolving the population toward a good solution. The algorithm begins with an initial random population. In EvA, this population is a set of vectors $\mathcal{S}_0 = \{\mathbf{s}_{i,0}\}_{i=1}^{n_p}$ with random elements, where n_p is the number of candidates in the population. By definition, our candidates always encode a valid perturbation—the budget of the perturbation is enforced by the length of each candidate vector.

In each iteration, candidates are evaluated using the fitness function. Based on the fitness scores, an elite sub-population of parents and new children is selected to proceed to the next iteration, while the rest of the population is removed. To create a new child, parents are selected through a tournament selection process: in each tournament, n_{tour} random parents are chosen, and the best among them is selected for crossover and subsequent mutation. This process repeats for t generations.

Sparse encoding of the attack. The population in our framework is an encoding of the perturbation matrix. The naive way for encoding this problem is to create a boolean vector of size N^2 encoding which entries are flipped, which results in memory complexity of $\mathcal{O}(|\mathcal{S}|N^2)$ where $|\mathcal{S}|$ is the population size. Instead, we introduce an approach that leverages the sparse nature of the solution and reduces the complexity to $\mathcal{O}(|\mathcal{S}| \cdot \epsilon \cdot |\mathcal{E}[\mathcal{V}_{\text{att}} : \mathcal{V}]|)$. In this encoding, instead of retaining all possible edges, we only keep the indices of the edges we want to flip. Therefore, any element in the population $\mathbf{z} \in \mathcal{S}$ is a vector of $p = \lfloor \epsilon \cdot |\mathcal{E}[\mathcal{V}_{\text{att}} : \mathcal{V}]| \rfloor$ dimensions where each entity of it is an index in adjacency matrix $\mathbf{z}[i] \in \{1, \dots, n(n-1)/2\}$ with $n = |\mathcal{V}|$. We use diagonal enumeration of an upper triangular $n \times n$ matrix as the encoding (see § B). The perturbation vector can contain repeated elements. During the evaluation of the vector, we transform it to a perturbation matrix $\mathbf{P}_{\mathbf{z}}$, and we compute the perturbed adjacency $\tilde{\mathbf{A}} = \mathbf{A} \oplus \mathbf{P}_{\mathbf{z}}$. All the mentioned computations are in sparse representation, and each individual of the population takes $\mathcal{O}(\delta)$ space. Moreover, with this encoding, we directly enforce the global budget since the size of each individual in the population is by design the number of allowed perturbations.

Acceptable fitness functions. The fitness function in GA is accessed in a black-box manner. Therefore, properties like differentiability are not a requirement, which allows us to use the accuracy directly. However, for scenarios like targeted attacks, where the objective is to only misclassify a single node, the 0-1 loss function is not a suitable fitness function. In other words, with the 0-1 loss, random search and GA are practically equivalent. Ideally, small changes in the solution should be reflected in the fitness function as well. This sensitivity to various individuals prevents GA from remaining in local optima. In § 5, we discuss the choice of fitness in targeted attacks.

Drawbacks. The aforementioned setup is the very baseline variant of EvA. While already effective (outperforming SOTA), in Fig. 2, we resolve several drawbacks by changing the definition of the initial population and the mutation function. In the baseline variant, a population is allowed to contain perturbations that are outside of the receptive field of the GNN for \mathcal{V}_{att} . This means that (at least for initial generations) a proportion of the attacking budget is wasted on ineffective perturbations. Even for perturbations connecting nodes with both ends outside of \mathcal{V}_{att} , the defender can easily revert them by memorizing the training subgraph. In § 4, we discuss further improvements.

4 ENHANCING THE SEARCH

In § 3 we defined the baseline evolutionary attack and discussed the possible drawbacks. As shown in Fig. 2 the baseline EvA already outperforms the SOTA. Additionally, with the following modifications we increase its effectiveness by a notable margin. The key insight is that the baseline attack, same as many gradient-based attacks defined their target space as the entire graph – entire space of $\frac{n}{2}(n-1)$ possible edges. As mentioned in § 2, perturbations that do have both endpoints in the training subgraph can be easily reverted just by memorizing the training subgraph. Additionally, perturbations outside of the receptive field of \mathcal{V}_{att} are a waste of budget as they do not affect the prediction of the target nodes.

Initial population. Our baseline initial population consists of random perturbations in the entire space of \mathcal{A} . This is a naive approach that disregards closeness to \mathcal{V}_{att} . Instead, we restrict the initial population to have at least one endpoint in \mathcal{V}_{att} . This can easily be done by randomly sampling both endpoints, one inside \mathcal{V}_{att} and one in \mathcal{V} , and then mapping the edges back to the indices via Π .

Targeted and adaptive mutation. After initialization, another way to balance the exploration and exploitation power of the algorithm is by introducing diversity in the population. In the baseline uniform mutation function, we change each edge to another random edge in \mathcal{V} with some mutation probability p . Same as in initialization, we define the “targeted mutation (TM)” by restricting the new mutated edge to have at least one end-point in \mathcal{V}_{att} . Remarkably, this modification shows a significant improvement as shown in Fig. 2. Furthermore, when the attack succeeds in altering a node’s prediction, additional perturbations connected to it do not gain any more performance. Therefore, we exclude them from the endpoint that was restricted to \mathcal{V}_{att} . Notably, we still allow those nodes to connect to other nodes in \mathcal{V}_{att} as they can also increase the misclassification risk for other nodes. We call the latter approach “adaptive targeted mutation” (ATM).

Stacking perturbations. Each population (at each iteration of EvA) needs to evaluate every individual. This means that each individual requires a forward pass on the perturbed graph. As mentioned before, our population takes $\mathcal{O}(\delta)$ memory, and during the evaluation, we still maintain the sparse representation of the graph. Therefore, if the memory budget allows, we can evaluate several perturbations at once by combining perturbed graphs into one large (disconnected) graph and running only one forward pass. In practice, for small datasets like CoraML, we only run one forward pass per iteration, as the entire population of 1024 individuals can be evaluated once.

Fitness Function. As we mentioned in § 1, one of the problems with gradient-based methods is finding a differentiable proxy aligned with the main objective. To further understand the effect of the loss function on attacks, we conducted an additional experiment where we replaced the fitness function of EvA with the cross-entropy and margin-based loss functions, which have become popular in adversarial attacks as surrogates for accuracy. This experiment seeks to evaluate the effect of the fitness function on attack performance. The results, shown in Fig. 2, indicate that cross-entropy does not use the budget effectively. On the contrary, the margin-based loss provides a well-correlated surrogate loss. Since PRBCD also uses the margin-based loss, we see that the main reason for the large gap to EvA is not the loss function. We hypothesize that EvA, leveraging the exploratory capabilities of genetic algorithms, can more effectively explore the solution space and avoid bad local optima, while PRBCD gets stuck.

Sensitivity. The fitness landscape should be sensitive – small changes in the solution should ideally result in (at least some) changes in the fitness score. For EvA, higher sensitivity results in a better selection of the population for breeding and distinguishes even the smallest advantage of a specific individual. We empirically show that accuracy has enough sensitivity for the global attack and low to medium size budget. However, as we discuss in § 5 for targeted attacks, the variability of the fitness function decreases to two values $\{0, 1\}$. We discuss further in § D.1 why low sensitivity of the objective function makes GA-based methods close to random search.

Effect of scaling. A larger population provides greater diversity among solutions, which helps prevent early convergence to sub-optimal solutions, therefore the population size has a considerable impact on the performance of EvA. To observe this effect, we conducted experiments by changing the population size while keeping other parameters fixed on the PubMed dataset. For a fair comparison, we also attempted to scale PRBCD by increasing the number of steps and the size of the block

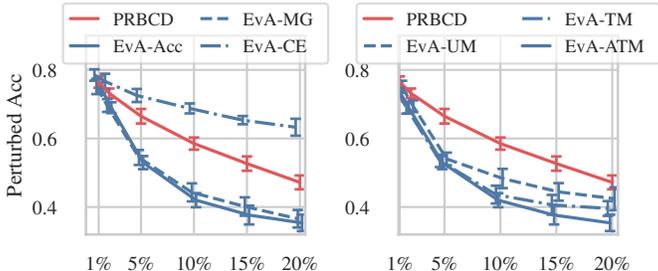


Figure 2: Effect of optimizing for different objective functions (left) and the influence of mutation type on EvA performance (right).

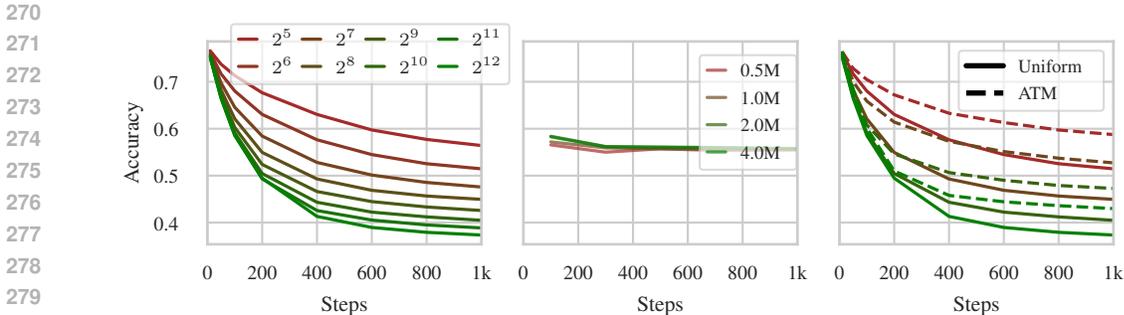


Figure 3: Effect of scaling on EvA and PRBCD performance (left, middle), the effect of mutation type and scaling on at 0.1% budget (right) on Pubmed dataset.

coordinate subspace. In this experiment, we exponentially increased the block size, starting from 500 up to 4 million. As shown in Fig. 3 (left), increasing the population size improves EvA’s ability to find better solutions by exploring the search space more effectively. Increasing the number of steps also increases the success rate of the attack. In contrast, PRBCD does not achieve further improvement by increasing the block size or the number of training steps.

We further investigate the effects of scaling and mutation types together. Fig. 3 (right) shows that adaptive targeted mutation can consistently enhance performance across all population sizes and outperform the uniform approach. This highlights the importance of selecting effective mutation strategies. Moreover, further exploration and refinement of mutation techniques could reveal more effective mutations, which could be explored in future studies.

5 OTHER OBJECTIVES

Local attacks. Gosch et al. (2024) argue that the perturbations within a global budget can still cause meaningful changes to the graph structure. For example, a perturbation might add edges to the node, increasing its degree to more than twice its current level while staying within the global budget. This can drastically alter the graph structure locally around that node, making the attack noticeable or, at the very least, impacting the graph’s structural semantics. Therefore, they argue for a threat model that, in addition to a global budget, has a local budget that limits the number of perturbations per node to an ϵ_{loc} proportion of its degree.

Local constrained mutation. We enforce this constraint as a new mutation applied before finalizing the population. In our mutation, we count the row (or column) summation of the perturbation matrix which quantifies the number of edges added to (or removed from) each node. Calling the nodes with perturbation degree higher than $\epsilon_{loc} \cdot \deg(v_i)$ (the local perturbation budget) as “violating”, we run an iterative refinement procedure where at each step we remove one edge from violating nodes and insert a non-violating edge instead. This refinement procedure continues until the local constraints are satisfied. Additionally, we rewrite the adaptive targeted mutation to account for the local budget – we restrict the mutation edges to those with remaining local budget for both end-points.

Targeted attacks. The targeted attack aims at one node to misclassify it with the least possible number of perturbations. With the discussion in § 4 the binary objective does not capture differences between different solutions. Therefore we use a proxy tanh-margin loss as the fitness function.

5.1 ATTACKING NOVEL OBJECTIVES

In cases where the objective is not differentiable (e.g. accuracy), to apply gradient-based attacks, we need to find a differentiable surrogate that approximates the original objective. This is already discussed in § 4. Using these attacks becomes even more challenging when the attack objective is complicated and defined through several non-differentiable components (e.g., quantile computation or majority voting). Since our method nullifies the need for information from gradients, we can easily

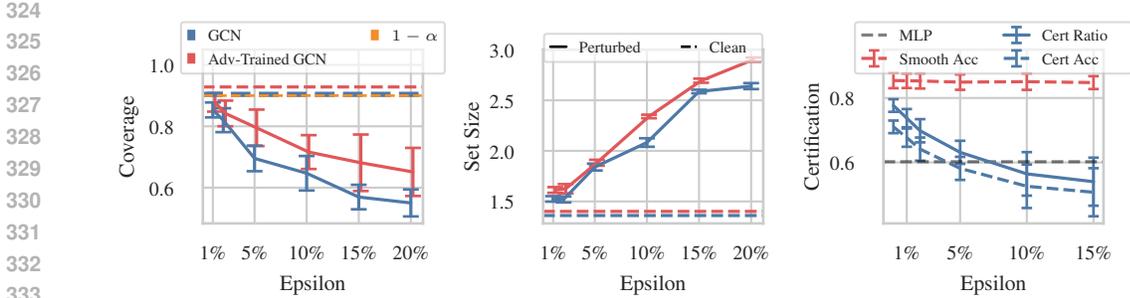


Figure 4: The conformal coverage (left) and conformal set size attack (middle) on Vanilla and adversarially trained GCN. The certificate attack (right) on GCN. All plots are for CoraML.

optimize for novel complex objectives. We define three new attacks on graphs: reducing the certified ratio of a smoothing-based model, decreasing the coverage, and the set size of conformal sets.

Attacking randomized smoothing-based certificates. A robustness certificate guarantees that the prediction of the classifier remains the same within the threat model. One way to obtain such a guarantee (for a black-box model) is through randomized smoothing. A smoothing scheme ξ is a random function mapping an input \mathbf{x} to a nearby point \mathbf{x}' (e.g. additive isotropic Gaussian noise $\mathbf{x}' = \xi(\mathbf{x}) = \mathbf{x} + \epsilon$, where $\epsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$). The convolution of the smoothing scheme and the classifier $\Pr[f(\mathbf{x} + \epsilon) = y]$ changes slowly around \mathbf{x} and this allows us to bound the worst-case minimum of the smooth prediction probability within \mathcal{B} . If this minimum is above 0.5, we can certify that the smooth model returns the same label for any $\tilde{\mathbf{x}} \in \mathcal{B}(\mathbf{x})$. A possible adversarial objective is to reduce the number of nodes that are certified (a.k.a. certified ratio).

For certifying a prediction we compute the majority vote – the probability that the classifier predicts the top class for randomized $\mathbf{x}' \sim \xi(\mathbf{x})$. Then, we find a lower bound for this probability within the perturbation ball (see § D). Exact computation of the majority vote is generally intractable. Instead, we use Monte-Carlo (MC) sampling. These operations are not directly differentiable.

A naive implementation of the certified ratio objective is to compute n_{mc} random samples for each candidate perturbation $\tilde{\mathbf{A}}$. This makes the attack extremely slow as in each iteration, we need $n_p \cdot n_{mc}$ samples, and for each Monte Carlo sample, we need n^2 samples from the Bernoulli distribution. Since statistical rigor is not crucial during the attack, we employ an efficient sampling strategy where we start with initial samples from clean \mathbf{A} , and for each perturbation, we only resample for the edges in $\tilde{\mathbf{A}} \triangle \mathbf{A}$. We use the stacked inference technique (see § 4) on MC samples which ultimately reduces the computation to one inference per each perturbation $\tilde{\mathbf{A}}$. Moreover, the certified radius is only a function of the smooth classifier’s probability and it is non-decreasing w.r.t. it. This allows us to reduce all certificate computations to one binary search for the minimum required probability. Then the objective is to minimize the number of nodes with probability above this threshold. We further discuss this attack in § D.

Attacking conformal prediction. Instead of label prediction, conformal prediction (CP) returns prediction sets that are guaranteed to include the true label with $1 - \alpha$ probability. This post-hoc statistical method treats the model as a black-box and requires only a calibration set of labeled points whose labels were not used during model training. CP is applicable in both inductive and transductive Graph Neural Networks (GNNs) under the assumption of node-exchangeability (Zargarbashi & Bojchevski, 2024). Adversarial attacks on conformal prediction aim to decrease the empirical coverage by perturbing the input. In addition, we also define an attack that reduces the applicability of the prediction sets by increasing the average set size.

To compute prediction sets we need to compute a quantile from the set of true calibration conformity scores and compare the scores of the test node to the quantile threshold. This operation is again not directly differentiable which is not a problem for EvA. In our experimental setup, the defender calibrates on a random subset of \mathcal{V}_u (besides the test, this is the only set with labels unseen by the model). Assuming that the unlabeled and test nodes are originally exchangeable (node-exchangeability), the conformal guarantee is valid in the inductive setup upon recalibration on the clean graph. By

378 perturbing the edge structure we can easily break this guarantee. Therefore our objective is to change
 379 the edge structure such that the coverage is minimized. Intuitively, this requires maximizing the
 380 distribution shift between the test and calibration scores. We can perform conformal prediction for
 381 each individual in the population, and we set the coverage of \mathcal{V}_{att} as the objective function.

382 Since we don't know the exact subset of the unlabeled nodes taken as calibration, we can use the
 383 unlabeled set entirely as the calibration set. Given that the defender will randomly sample from
 384 unlabeled nodes during the calibration, the coverage remains roughly the same for exchangeable
 385 subsets of \mathcal{V}_u (Berti & Rigo, 1997). To the best of our knowledge, so far this is the only adversarial
 386 attack on the graph structure to break conformal inductive GNNs. Similarly, by changing the objective
 387 to the negative average set size, we can easily attack the usability of prediction sets (see Fig. 4).
 388

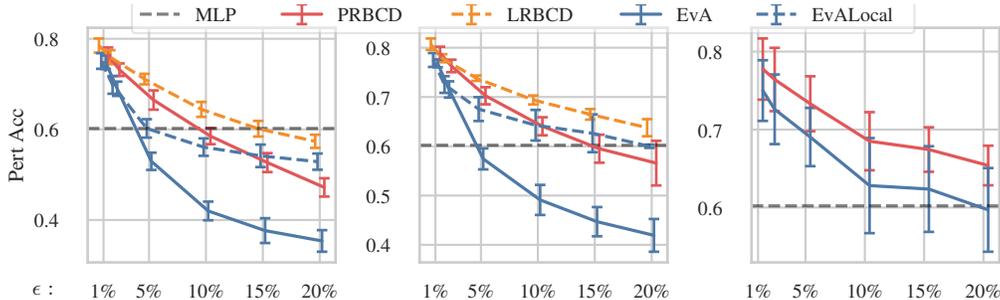
389 6 EMPIRICAL RESULTS

390 With our empirical evaluations (i) we show that current gradient-based attacks are still very far from
 391 optimal since EvA outperforms them by a notable margin. (ii) We show that EvA inherently results
 392 in attacks that perturb each node with less change in nodes' degree. This is even without posing local
 393 budget restrictions. (iii) We also show even by adding local restrictions EvA still outperforms other
 394 gradient-based local attacks. (iv) The effectiveness of EvA is consistent across various models, and
 395 vanilla or robust training setups. (v) With the black-box nature of the attack we introduce the first
 396 attack that reduces the certified ratio and the first attack that breaks conformal sets on graphs.
 397

398 **Experimental setup.** We evaluate EvA on common graph datasets: Cora-ML (McCallum et al.,
 399 2004), Citeseer (Sen et al., 2008), and PubMed (Namata et al., 2012). Shchur et al. (2018) show that
 400 GNN evaluation is sensitive to the initial train/val/test split. Therefore, we averaged our results for
 401 each dataset/model over five different data splits. In contrast with common GNN attacks, Gosch et al.
 402 (2024) show that transductive setup carries a false sense of robustness. In other words, trivially one
 403 can gain perfect robustness just by memorizing the clean data; models with robust and self-training
 404 also show to exploit this flaw. Following them, we report our results in an inductive setting. We
 405 divide graph nodes into four subsets: training, validation, and testing, each with 10% of the nodes
 406 and we leave the remaining 60% as unlabeled data. For completeness, in § A we compare attacks in
 407 the transductive setting as well, where again EvA is more effective.

408 Following Lingam et al. (2023), we maintain the distribution of labels for sampling train, validation,
 409 and test nodes. This provides a more realistic scenario compared to commonly used methods, such as
 410 sampling for training and validation with the same count probability for each class. For completeness
 411 in § A we report various sampling setups. However the this does not change the order between
 412 methods. Further information about the model and hyperparameters can be found in § C.

413 **Attacking vanilla models.** As shown in Fig. 5, EvA outperforms the SOTA attack PRBCD by a
 414 significant margin. This comparison remains consistent across various datasets and models. We report
 415 these results extensively in § A. Interestingly, we show that in many vanilla and robust models, a very
 416 small budget $\epsilon \sim 0.05$ EvA drops the accuracy below the level of the MLP model. This is a condition
 417 where the model leveraging the structure works worse than a model that completely ignores edges.
 418



429
 430 Figure 5: Performance of EvA on CoraML. From left to right the results are on Vanilla GCN,
 431 adversarially trained GCN using PRBCD, and Soft-Median-GDC model.

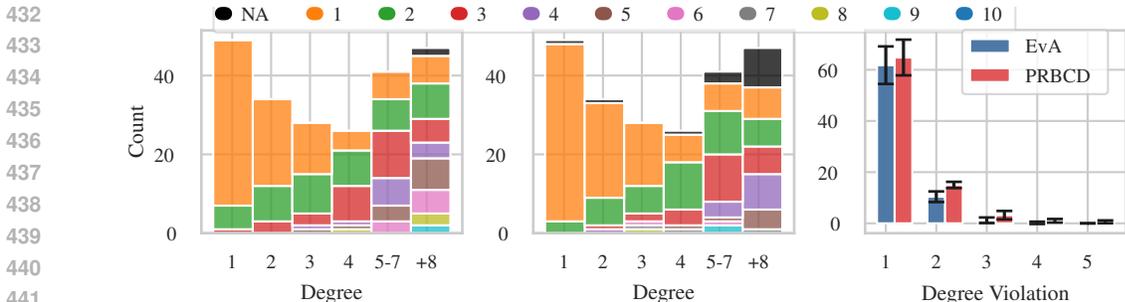


Figure 6: The number of perturbations (in different colors) that have been used by EvA (left) and PRBCD (middle) to target nodes with a specific degree. The right figures present the number of violations that EvA, PRBCD introduce (for $\epsilon_{loc} = 0.5$). NA (black) indicates a failed attack.

The SoftMedian model seems to show an inherent robustness to both EvA and PRBCD. Therefore, to break the model below the accuracy of MLP, we require ≥ 0.2 perturbation budget. Even in the SoftMedian model, our attack is significantly more effective in comparison to PRBCD.

Adversarially trained models. Similar to vanilla models, EvA outperforms other approaches for robust models by a notable margin. As expected, models trained with EvA were shown to be more robust, and other attacks were less effective to them. However, this additional robustness is not significant. Table 9 (§ A) compares attacks in models with different adversarial training.

Local attacks. Building upon the discussion in § 5, we enforce the local constraint by adding a mutation function that iteratively removes edges exceeding the local budget. The local variant of EvA also shows to be consistently better than the local LRBCD attack. For the PubMed dataset EvA’s effectiveness has a slower trend by increasing the budget ϵ . On PubMed (see § A), when enforcing locally constrained mutation, the number of reconsidered edges increases due to the density of the graph. This makes the search significantly harder. On other datasets, though, EvA shows consistently better results and, more importantly, sharper decrease at lower ϵ .

Targeted attack. We perform attacks on each node separately, with varying budgets from one to a maximum of 10 edges, until the prediction changes. As we discussed in § 5, the accuracy on one node is non-expressive. Therefore we use the tanh-Margin proxy loss. Fig. 6 (left and middle) compares EvA and PRBCD in targeted attack. Our results show that PRBCD performs better with a budget of one, but is outperformed by EvA for budgets of two and higher. For instance, on the CoraML dataset PRBCD fails to modify 16 nodes with a maximum of 10 changes (NA, black), whereas this number is reduced to only 2 nodes for EvA. This result is expected due to the combinatorial nature of the problem: for budgets up to two, a greedy approach can find the optimal solution, but as the budget increases beyond three, the problem becomes significantly more complex.

Attacking certificates. As shown in Fig. 4 (right) we reduce the certified ratio - the number of nodes that the certificate can guarantee for the specified threat model - to a ratio below the accuracy of the MLP model. The MLP model here is a baseline as it is robust to any structure perturbation by trivially ignoring edges. We report the ratio certified by sparse smoothing (Bojchevski et al., 2020) with $p_- = 0.4$, and $p_+ = 2 \times 10^{-5}$. Here p_+ , and p_- are Bernoulli parameters of flipping a zero or one. We reported the result for $\mathcal{B}_{0,3}$ which means 0 edge addition and 3 deletions. While we aim to decrease the certified ratio, a direct outcome is that the certified accuracy drops. For a 5% budget, the certified accuracy drops below MLP, which is resilient to any structural perturbation by definition. Notably, the clean (smooth) accuracy stays the same making this attack less noticeable. For this experiment we used the GPRGNN model with robust training using PRBCD attack.

Attacking conformal prediction. We report the first structure attack to inductive conformal GNN (Zargarbashi & Bojchevski, 2024). As shown in Fig. 4 (right) the coverage drops quickly as we increase the perturbation budget. As expected, in an adversarially trained model, we observe a slower decrease in the empirical coverage. Another interesting objective to attack is increasing the set size since it affects the usability of the prediction set. In Fig. 4 (middle) we show that both vanilla and robust models are vulnerable to this attack.

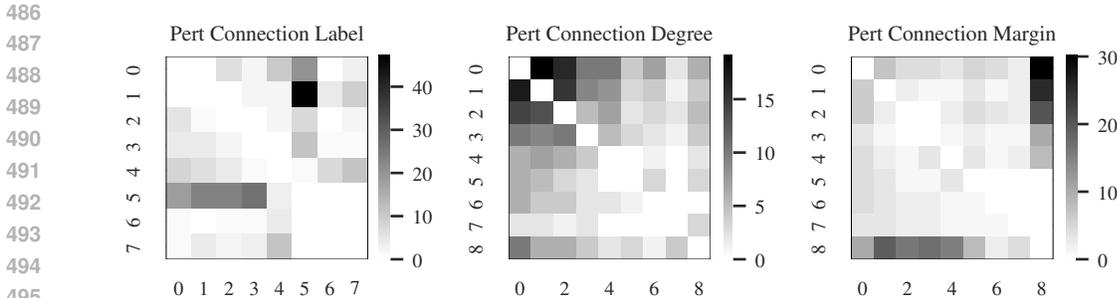


Figure 7: The upper triangle of each heatmap represents the perturbation connections for PRBCD, the lower triangle corresponds to the same for EvA, and the diagonal is set to zero.

Local degree violation. In this experiment, we did not enforce the local degree constraint in both EvA and PRBCD. However, we compared the final solutions to assess how often the solutions violated this constraint. Fig. 6 (right) shows that EvA generally produces more diverse attacks, utilizing the global budget more efficiently, which leads to improved performance. Overall, EvA violates less for any degree number and therefore in total.

Label diversity. We further conduct an ablation study on the solutions found by EvA and PRBCD under a specific budget of 10%. In this experiment, we keep all hyperparameters of EvA and PRBCD fixed and run them across 10 different seeds. We then compare the average solutions generated by each adversary. The left figure in Fig. 7 shows the number of connections across different labels. In both cases, the methods focus more on label 5 than on the others, but EvA distributes the connections more uniformly compared to PRBCD. The middle figure illustrates the nodes with original degrees ranging from 1 to greater than 8. The results indicate that, in both attacks, most of the budget is spent connecting to low-degree nodes. However, compared to PRBCD, EvA allocates more of the budget to higher-degree nodes. Additionally, we calculate the margin loss for each node in the original graph and discretize them into eight levels. As shown in the right figure of Fig. 7, EvA allocates more of the budget to higher-margin nodes, resulting in a non-trivial solution that achieves a better optimum. Finally, it seems that EvA identifies solutions that differ from greedy-based heuristic, which usually only targets low-degree or low-margin nodes.

7 CONCLUSION

In contrast to gradient-based adversarial attacks on graph structure, we developed a new attack (EvA) based on a heuristic genetic algorithm. By eliminating differentiation, we can directly optimize for the objective of the adversary (e.g. the model’s accuracy). This black-box nature enables us to define complex adversarial goals, including attacks on robustness certificates and conformal prediction. Our novel attacks decrease the certified ratio, and conformal coverage, and increase the conformal set size. We propose an encoding that reduces the memory complexity of the attack to the same order as the perturbation budget which allows us to adapt to various computational constraints. Given the drastic decrease in the model’s accuracy by applying EvA, we highlight that even SOTA gradient-based attacks are far from optimal. Our main message is that search-based attacks are underexplored yet powerful as shown by our results.

Limitations. We use an off-the-shelf genetic algorithm. Surely, there is room for designing search algorithms specific to the domain of the problem or hybrids of gradient and evolutionary search. As the graph size increases, the search space expands exponentially which makes convergence harder. While we remove the white-box assumption, we still assume the adversary has full knowledge of the graph and labels (same as most other attacks). This limitation can be easily addressed in future work. EvA uses many forward passes through the model which can be unrealistic in some attack scenarios. We leave the design of a further query-efficient variant for the future.

540 ETHICS STATEMENT

541

542 In this paper, we propose an adversarial attack without white-box access. However, our main focus is
543 to point out the vulnerability of GNN models, opening the discussion on the need for more robust
544 and reliable models, our results can be used to exploit the vulnerability of current GNNs.

545

546 REPRODUCIBILITY

547

548 For reproducibility, we have uploaded the complete anonymized codebase on OpenReview.

549

550 REFERENCES

551

552 Anish Athalye, Nicholas Carlini, and David Wagner. Obfuscated gradients give a false sense of
553 security: Circumventing defenses to adversarial examples. In *International conference on machine
554 learning*, pp. 274–283. PMLR, 2018.

555

556 Patrizia Berti and Pietro Rigo. A glivenko-cantelli theorem for exchangeable random variables.
557 *Statistics & probability letters*, 32(4):385–391, 1997.

558

559 Aleksandar Bojchevski, Johannes Gasteiger, and Stephan Günnemann. Efficient robustness certifi-
560 cates for discrete data: Sparsity-aware randomized smoothing for graphs, images and more. In
International Conference on Machine Learning, pp. 1003–1013. PMLR, 2020.

561

562 Eli Chien, Jianhao Peng, Pan Li, and Olgica Milenkovic. Adaptive universal generalized pagerank
563 graph neural network, 2021.

564

565 Hanjun Dai, Hui Li, Tian Tian, Xin Huang, Lin Wang, Jun Zhu, and Le Song. Adversarial attack on
566 graph structured data. In *International conference on machine learning*, pp. 1115–1124. PMLR,
2018.

567

568 Simon Geisler, Tobias Schmidt, Hakan Şirin, Daniel Zügner, Aleksandar Bojchevski, and Stephan
569 Günnemann. Robustness of graph neural networks at scale. *Advances in Neural Information
570 Processing Systems*, 34:7637–7649, 2021.

571

572 Simon Geisler, Tobias Schmidt, Hakan Şirin, Daniel Zügner, Aleksandar Bojchevski, and Stephan
573 Günnemann. Robustness of graph neural networks at scale, 2023. URL [https://arxiv.org/
574 abs/2110.14038](https://arxiv.org/abs/2110.14038).

575

576 Lukas Gosch, Simon Geisler, Daniel Sturm, Bertrand Charpentier, Daniel Zügner, and Stephan
577 Günnemann. Adversarial training for graph neural networks: Pitfalls, solutions, and new directions.
In *37th Conference on Neural Information Processing Systems (Neurips)*, 2023.

578

579 Lukas Gosch, Simon Geisler, Daniel Sturm, Bertrand Charpentier, Daniel Zügner, and Stephan
580 Günnemann. Adversarial training for graph neural networks: Pitfalls, solutions, and new directions.
Advances in Neural Information Processing Systems, 36, 2024.

581

582 John H Holland. Genetic algorithms and adaptation. *Adaptive control of ill-defined systems*, pp.
583 317–333, 1984.

584

585 Vijay Lingam, Mohammad Sadeh Akhondzadeh, and Aleksandar Bojchevski. Rethinking label
586 poisoning for gnns: Pitfalls and attacks. In *The Twelfth International Conference on Learning
587 Representations*, 2023.

588

589 Andrew McCallum, Kamal Nigam, Jason D. M. Rennie, and Kristie Seymore. Automating the
construction of internet portals with machine learning. *Information Retrieval*, 3:127–163, 2004.

590

591 Galileo Namata, Ben London, Lise Getoor, and Bert Huang. Query-driven active surveying for
collective classification. 2012.

592

593 Prithviraj Sen, Galileo Namata, Mustafa Bilgic, Lise Getoor, Brian Gallagher, and Tina Eliassi-Rad.
Collective classification in network data. 2008.

- 594 Oleksandr Shchur, Maximilian Mumme, Aleksandar Bojchevski, and Stephan Günnemann. Pitfalls
595 of graph neural network evaluation. *arXiv preprint arXiv:1811.05868*, 2018.
596
- 597 Lichao Sun, Yingtong Dou, Carl Yang, Kai Zhang, Ji Wang, Philip S. Yu, Lifang He, and Bo Li.
598 Adversarial attack and defense on graph data: A survey. *IEEE Transactions on Knowledge and
599 Data Engineering*, 35(8):7693–7711, 2023. doi: 10.1109/TKDE.2022.3201243.
- 600 Nihat Engin Toklu, Timothy Atkinson, Vojtěch Micka, Paweł Liskowski, and Rupesh Kumar Srivas-
601 tava. Evotorch: scalable evolutionary computation in python. *arXiv preprint arXiv:2302.12600*,
602 2023.
603
- 604 Yexin Wang, Zhi Yang, Junqi Liu, Wentao Zhang, and Bin Cui. Scapin: Scalable graph structure
605 perturbation by augmented influence maximization. *Proc. ACM Manag. Data*, 1(2), June 2023.
606 doi: 10.1145/3589291. URL <https://doi.org/10.1145/3589291>.
- 607 Kaidi Xu, Hongge Chen, Sijia Liu, Pin-Yu Chen, Tsui-Wei Weng, Mingyi Hong, and Xue Lin.
608 Topology attack and defense for graph neural networks: An optimization perspective. *arXiv
609 preprint arXiv:1906.04214*, 2019.
- 610 Soroush H Zargarbashi and Aleksandar Bojchevski. Conformal inductive graph neural networks.
611 *arXiv preprint arXiv:2407.09173*, 2024.
612
- 613 Chenhan Zhang, Shiyao Zhang, James J. Q. Yu, and Shui Yu. Sam: Query-efficient adversarial
614 attacks against graph neural networks. *ACM Trans. Priv. Secur.*, 26(4), November 2023. ISSN
615 2471-2566. doi: 10.1145/3611307. URL <https://doi.org/10.1145/3611307>.
- 616 Jianfu Zhang, Yan Hong, Dawei Cheng, Liqing Zhang, and Qibin Zhao. Hierarchical attacks on large-
617 scale graph neural networks. In *ICASSP 2024 - 2024 IEEE International Conference on Acoustics,
618 Speech and Signal Processing (ICASSP)*, pp. 7635–7639, 2024. doi: 10.1109/ICASSP48485.2024.
619 10448076.
620
- 621 Guanghui Zhu, Mengyu Chen, Chunfeng Yuan, and Yihua Huang. Simple and efficient partial
622 graph adversarial attack: A new perspective, 2023. URL [https://arxiv.org/abs/2308.
623 07834](https://arxiv.org/abs/2308.07834).
- 624 Daniel Zügner, Amir Akbarnejad, and Stephan Günnemann. Adversarial attacks on neural networks
625 for graph data. In *Proceedings of the 24th ACM SIGKDD international conference on knowledge
626 discovery & data mining*, pp. 2847–2856, 2018.
627
- 628 Daniel Zügner, Oliver Borchert, Amir Akbarnejad, and Stephan Günnemann. Adversarial attacks
629 on graph neural networks: Perturbations and their patterns. *ACM Transactions on Knowledge
630 Discovery from Data (TKDD)*, 14(5):1–31, 2020.
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647

A SUPPLEMENTARY EXPERIMENTS

Transductive Setting. In § 6, we discussed that evaluating robustness in transductive setup is flawed since trivial robustness can be gained just by memorization of the clean graph (Gosch et al., 2024). However, for completeness, Table 1 reports the attacks’ effectiveness in this setup. Consistent with other experiments, here also EvA outperforms SOTA. We provide the result for EvA, PRBCD, and for completeness, we also provide the result for PGA, which is a more recent attack.

Table 1: Classification accuracy (%) on the CoraML dataset in the transductive setting under different attacks and perturbation levels ϵ . Results are averaged over multiple runs with standard deviations.

Model	Attack	ϵ					
		0.01	0.02	0.05	0.10	0.15	0.20
GCN	LRBCD	80.0 \pm 2.0	78.0 \pm 2.0	73.0 \pm 1.0	66.0 \pm 1.0	61.0 \pm 2.0	57.0 \pm 2.0
	PRBCD	79.0 \pm 2.0	77.0 \pm 2.0	72.0 \pm 2.0	65.0 \pm 2.0	60.0 \pm 2.0	56.0 \pm 2.0
	EvA	77.0 \pm 2.0	74.0 \pm 2.0	66.0 \pm 2.0	60.0 \pm 2.0	59.0 \pm 3.0	57.0 \pm 3.0
GPRGNN	LRBCD	79.0 \pm 3.0	77.0 \pm 3.0	72.0 \pm 4.0	63.0 \pm 7.0	55.0 \pm 12.0	49.0 \pm 16.0
	PRBCD	79.0 \pm 3.0	76.0 \pm 4.0	70.0 \pm 5.0	62.0 \pm 7.0	55.0 \pm 10.0	50.0 \pm 13.0
	EvA	77.0 \pm 3.0	73.0 \pm 5.0	64.0 \pm 6.0	57.0 \pm 10.0	53.0 \pm 13.0	50.0 \pm 16.0

Inductive Setting. Here, we present additional results specifically for the inductive setting. Unlike the transductive setup, where robustness can be misleadingly achieved through memorization of the clean graph, the inductive framework provides a more comprehensive assessment of model performance in real-world scenarios. In this section, we detail the effectiveness of our method compared to other approaches.

Table 2: Classification accuracy (%) on the CoraML dataset in the inductive setting under different attacks and perturbation levels ϵ . Results are averaged over multiple runs with standard deviations.

Model	Attack	ϵ					
		0.01	0.02	0.05	0.10	0.15	0.20
APPNP	EvA	76.65 \pm 1.32	71.03 \pm 1.44	56.51 \pm 1.60	49.32 \pm 1.84	44.77 \pm 2.04	41.42 \pm 1.41
	PRBCD	78.65 \pm 0.99	75.30 \pm 1.27	68.75 \pm 1.22	61.57 \pm 1.65	55.44 \pm 1.58	49.96 \pm 2.42
GAT	EvA	64.20 \pm 1.89	58.51 \pm 2.45	40.99 \pm 1.60	15.30 \pm 4.47	9.40 \pm 6.83	8.11 \pm 6.65
	PRBCD	70.07 \pm 2.82	66.55 \pm 2.21	58.58 \pm 3.33	49.61 \pm 6.55	39.86 \pm 6.78	36.94 \pm 7.09
GCN	EvA	74.80 \pm 1.50	68.97 \pm 1.58	52.95 \pm 1.91	41.99 \pm 2.06	37.65 \pm 2.74	35.37 \pm 2.38
	PRBCD	76.44 \pm 1.64	73.17 \pm 1.39	66.48 \pm 2.13	58.51 \pm 1.77	52.67 \pm 2.09	47.19 \pm 2.02
GPRGNN	EvA	72.53 \pm 4.11	66.83 \pm 4.54	51.53 \pm 5.57	42.21 \pm 8.52	37.01 \pm 9.83	34.52 \pm 9.83
	PRBCD	74.95 \pm 3.08	71.67 \pm 2.76	64.84 \pm 4.18	57.94 \pm 4.55	53.24 \pm 5.20	48.68 \pm 6.52

Stratified sampling. Although unrealistic, in Table 8 we compare attacks in case the models are trained train/val/test sampled with the same number of nodes across different classes. Consistent with other results, EvA shows to be better here as well.

Attacking accuracy of vanilla and robust models. Table 9 compares EvA with SOTA PRBCD, and LRBCD. We compare both three attacks on vanilla models or models trained with adversarial examples of either of the attacks. Across all setups, EvA shows a comparably better performance.

A.1 EXPERIMENTS

Vanilla Models Experiments. For the experimental results, we mainly focus on the inductive setting introduced by (Lingam et al., 2023), where during training, we only use \mathcal{G}_{tr} , and during the attack, we target \mathcal{G}_{test} . We also provide results for the transductive setting, showcasing that our attack outperforms previous gradient-based methods, independent of the training setting. We conduct

Table 3: Classification accuracy (%) on the CoraML dataset in the inductive setting under different attacks and perturbation levels ϵ . Results are averaged over multiple runs with standard deviations.

Model	Attack	ϵ					
		0.01	0.02	0.05	0.10	0.15	0.20
GCN	EvA	74.80 \pm 1.50	68.97 \pm 1.58	52.95 \pm 1.91	41.99 \pm 2.06	37.65 \pm 2.74	35.37 \pm 2.38
	EvaLocal	75.09 \pm 1.73	69.82 \pm 1.96	60.21 \pm 2.04	56.09 \pm 1.93	54.16 \pm 2.48	52.88 \pm 1.79
	LRBCD	78.51 \pm 1.56	75.94 \pm 1.54	71.10 \pm 1.16	64.41 \pm 1.65	60.14 \pm 1.73	57.37 \pm 1.45
	PRBCD	76.44 \pm 1.64	73.17 \pm 1.39	66.48 \pm 2.13	58.51 \pm 1.77	52.67 \pm 2.09	47.19 \pm 2.02
	PGA	79.58 \pm 1.61	76.92 \pm 1.73	70.94 \pm 1.89	64.62 \pm 1.92	60.46 \pm 2.25	57.54 \pm 2.46
GPRGNN	EvA	72.53 \pm 4.11	66.83 \pm 4.54	51.53 \pm 5.57	42.21 \pm 8.52	37.01 \pm 9.83	34.52 \pm 9.83
	EvaLocal	73.31 \pm 3.30	67.26 \pm 4.17	58.29 \pm 7.96	53.38 \pm 11.42	51.10 \pm 12.66	49.96 \pm 13.63
	LRBCD	77.51 \pm 1.81	74.80 \pm 1.41	68.83 \pm 1.90	62.56 \pm 1.71	59.07 \pm 1.53	55.66 \pm 1.71
	PRBCD	74.95 \pm 3.08	71.67 \pm 2.76	64.84 \pm 4.18	57.94 \pm 4.55	53.24 \pm 5.20	48.68 \pm 6.52
	PGA	78.55 \pm 3.03	75.33 \pm 3.69	68.63 \pm 5.11	61.55 \pm 6.97	56.60 \pm 8.52	54.91 \pm 7.46

Table 4: Classification accuracy (%) on the Citeseer dataset in the inductive setting under different attacks and perturbation levels ϵ . Results are averaged over multiple runs with standard deviations.

Model	Attack	ϵ					
		0.01	0.02	0.05	0.10	0.15	0.20
APPNP	EvA	-	-	74.29 \pm 0.88	65.00 \pm 1.15	59.76 \pm 2.33	54.76 \pm 1.19
	PRBCD	87.26 \pm 0.90	85.48 \pm 1.49	81.79 \pm 1.08	76.55 \pm 0.68	72.44 \pm 1.66	69.29 \pm 1.69
GAT	EvA	-	-	67.14 \pm 3.65	51.19 \pm 4.21	37.74 \pm 4.94	27.62 \pm 9.49
	PRBCD	84.52 \pm 2.27	82.62 \pm 2.20	76.55 \pm 5.59	70.00 \pm 6.09	67.02 \pm 4.27	63.15 \pm 4.19
GCN	EvA	86.67 \pm 1.71	82.86 \pm 2.12	72.74 \pm 2.74	58.33 \pm 3.01	49.76 \pm 3.22	44.29 \pm 3.33
	PRBCD	87.38 \pm 1.81	85.83 \pm 2.43	80.95 \pm 2.06	74.29 \pm 4.22	69.76 \pm 4.34	67.62 \pm 4.96
GPRGNN	EvA	87.26 \pm 2.75	83.81 \pm 2.50	73.45 \pm 3.17	61.43 \pm 4.66	55.48 \pm 3.84	50.12 \pm 4.86
	PRBCD	88.45 \pm 2.29	86.31 \pm 2.45	82.02 \pm 2.61	77.14 \pm 2.84	73.93 \pm 3.89	69.64 \pm 3.47

experiments on the Cora-ML Citeseer, and Pubmed datasets, trained GCN, GPRGNN, APPNP and GAT. We run the attack for six different budgets (0.01, 0.02, 0.05, 0.1, 0.15, 0.2). Further details on training and attack hyperparameters are provided in C. We also use EvoTorch (Toklu et al., 2023) to implement EvA.

Robust Models. Similarly, we provide the results for the adversarially trained model. In this case, during training, we use an adversarial attack at each step to attack \mathcal{G}_{tr} , and then we retrain the model on the adversarially perturbed graph $\tilde{\mathcal{G}}_{tr}$. The robust budget (ϵ_{robust}) for adversarial attack during training was 0.02. This process repeats in each epoch of training until the model converges. We similarly use the inductive setting since, as Gosch et al. (2023) shows, in the transductive setup, the evaluation is flawed by a false sense of robustness. This originates from the fact that if, during the training process, the defender has access to perfect knowledge of all nodes in the graph, it can achieve perfect robustness by memorizing the clean structure of the graph in the model’s weights. Table 9 presents our results in this setting. As the results indicate, EvA outperforms all previous attacks, even in adversarially trained models.

Attacking Certificate. As we discussed, EvA still be used in the case that we don’t have access to a non-differentiable objective. In this experiment, we introduced a certificate attack which, to the best of our knowledge, is the first attack on certificates that reduces the certificate guarantees

Any possible appendices should be placed after bibliographies. If your paper has appendices, please submit the appendices together with the main body of the paper. There will be no separate supplementary material submission. The main text should be self-contained; reviewers are not obliged to look at the appendices when writing their review comments.

Table 5: Classification accuracy (%) on the Citeseer dataset in the inductive setting under different attacks and perturbation levels ϵ . Results are averaged over multiple runs with standard deviations.

Model	Attack	ϵ					
		0.01	0.02	0.05	0.10	0.15	0.20
GCN	EvA	86.67 \pm 1.71	82.86 \pm 2.12	72.74 \pm 2.74	58.33 \pm 3.01	49.76 \pm 3.22	44.29 \pm 3.33
	EvaLocal	87.38 \pm 1.65	83.57 \pm 2.17	78.21 \pm 3.17	76.43 \pm 2.62	75.00 \pm 3.23	74.52 \pm 3.16
	LRBCD	88.45 \pm 2.17	86.43 \pm 2.71	83.69 \pm 2.48	80.12 \pm 3.30	78.45 \pm 3.89	75.36 \pm 4.81
	PRBCD	87.38 \pm 1.81	85.83 \pm 2.43	80.95 \pm 2.06	74.29 \pm 4.22	69.76 \pm 4.34	67.62 \pm 4.96
GPRGNN	EvA	87.26 \pm 2.75	83.81 \pm 2.50	73.45 \pm 3.17	61.43 \pm 4.66	55.48 \pm 3.84	50.12 \pm 4.86
	EvaLocal	87.50 \pm 2.27	84.29 \pm 2.04	80.48 \pm 3.96	77.86 \pm 4.56	76.43 \pm 5.06	75.12 \pm 6.57
	LRBCD	89.76 \pm 2.50	87.98 \pm 2.48	85.12 \pm 2.76	81.90 \pm 2.83	79.64 \pm 4.08	78.45 \pm 4.92
	PRBCD	88.45 \pm 2.29	86.31 \pm 2.45	82.02 \pm 2.61	77.14 \pm 2.84	73.93 \pm 3.89	69.64 \pm 3.47

Table 6: Classification accuracy (%) on the PubMed dataset in the inductive setting under different attacks and perturbation levels ϵ . Results are averaged over multiple runs with standard deviations.

Model	Attack	ϵ					
		0.01	0.02	0.05	0.10	0.15	0.20
APPNP	EvA	73.85 \pm 2.35	69.64 \pm 2.16	57.07 \pm 2.32	47.03 \pm 2.18	43.94 \pm 1.83	41.93 \pm 2.18
	PRBCD	75.54 \pm 2.34	72.44 \pm 2.28	65.14 \pm 2.26	57.15 \pm 2.59	51.04 \pm 2.79	45.75 \pm 2.60
GAT	EvA	69.15 \pm 1.83	64.62 \pm 1.81	52.17 \pm 1.71	33.62 \pm 2.05	26.62 \pm 3.74	24.21 \pm 4.27
	PRBCD	71.33 \pm 1.53	67.78 \pm 1.79	59.73 \pm 2.10	49.87 \pm 1.36	42.04 \pm 1.57	35.94 \pm 1.66
GCN	EvA	72.60 \pm 2.19	68.35 \pm 2.41	56.15 \pm 1.92	42.93 \pm 2.64	40.46 \pm 2.76	39.11 \pm 2.98
	PRBCD	74.99 \pm 1.99	71.90 \pm 2.03	64.16 \pm 2.32	55.54 \pm 2.79	49.32 \pm 2.66	43.90 \pm 3.09
GPRGNN	EvA	72.01 \pm 4.18	67.61 \pm 4.28	55.95 \pm 4.32	—	—	42.39 \pm 9.63
	PRBCD	74.37 \pm 3.40	71.66 \pm 3.55	64.51 \pm 4.94	56.21 \pm 6.46	50.26 \pm 7.41	45.81 \pm 8.47

B TECHNICAL DETAILS OF EVA

Mapping function: enumeration over \mathbf{A} For enumerating over \mathbf{A} , instead of using the row and column indices of the node to select, we introduced indexing. For a directed graph, the indexing starts from 0 to $n^2 - 1$. However, in an undirected graph, we only need the upper triangular part of the matrix \mathbf{A} . To achieve this, we use the following algebraic solution to find the row and column of the perturbation by referencing only the upper triangular indexing.

$$r = n - 2 - \left\lfloor \frac{\sqrt{-8l + 4n(n-1) - 7}}{2} - 0.5 \right\rfloor$$

$$c = 1 + l + r - \frac{n(n-1)}{2} + \left\lfloor \frac{(n-r)(n-r-1)}{2} \right\rfloor$$
(2)

The advantage of this solution is that it can also be implemented in a vectorized way, making everything parallelizable.

C DATASETS AND MODELS, AND HYPERPARAMETERS

C.1 STATISTICS OF DATASETS

In our experiments, we mainly conduct experiments on the commonly used graph datasets: Cora-ML, Citeseer, and PubMed, which are all representative academic citation networks. Their specific characteristics are as follows:

Cora-ML. The Cora-ML dataset contains 2,810 papers as nodes, with citation relationships between them as edges, resulting in 7,981 edges. Each paper is categorized into one of 7 classes corresponding

Table 7: Classification accuracy (%) on the PubMed dataset in the inductive setting under different attacks and perturbation levels ϵ . Results are averaged over multiple runs with standard deviations.

Model	Attack	ϵ					
		0.01	0.02	0.05	0.10	0.15	0.20
GCN	EvA	72.60 \pm 2.18	68.35 \pm 2.41	56.15 \pm 1.92	42.93 \pm 2.64	40.46 \pm 2.76	39.11 \pm 2.98
	EvaLocal	74.49 \pm 2.05	70.53 \pm 2.10	66.97 \pm 2.86	74.75 \pm 3.14	74.04 \pm 2.09	72.99 \pm 2.09
	LRBCD	74.89 \pm 2.04	71.48 \pm 2.49	65.68 \pm 2.90	60.24 \pm 3.15	56.81 \pm 3.02	54.07 \pm 2.99
	PRBCD	74.99 \pm 1.99	71.90 \pm 2.03	64.16 \pm 2.32	55.54 \pm 2.79	49.32 \pm 2.66	43.90 \pm 3.09
GPRGNN	EvA	72.01 \pm 4.18	67.61 \pm 4.28	55.95 \pm 4.32	—	—	42.39 \pm 9.63
	EvaLocal	73.36 \pm 3.71	69.68 \pm 3.93	65.61 \pm 6.75	70.64 \pm 2.85	72.27 \pm 5.06	71.40 \pm 5.41
	LRBCD	74.50 \pm 3.66	71.57 \pm 4.10	65.88 \pm 6.12	60.33 \pm 5.70	56.75 \pm 7.74	53.75 \pm 8.18
	PRBCD	74.37 \pm 3.40	71.66 \pm 3.55	64.51 \pm 4.94	56.21 \pm 6.46	50.26 \pm 7.41	45.81 \pm 8.47

Table 8: Classification accuracy (%) on the CoraML dataset in the inductive setting under different attacks and perturbation levels ϵ . Results are averaged over multiple runs with standard deviations. Training, validation, and test sets are stratified.

Model	Attack	ϵ					
		0.01	0.02	0.05	0.10	0.15	0.20
GCN	LRBCD	80.0 \pm 2.7	77.4 \pm 2.6	71.4 \pm 2.7	65.5 \pm 3.9	61.2 \pm 4.2	57.6 \pm 4.5
	PRBCD	78.7 \pm 2.8	75.3 \pm 3.4	67.9 \pm 3.4	59.5 \pm 3.5	53.0 \pm 4.1	48.4 \pm 3.7
	EvA	77.0 \pm 2.9	71.4 \pm 3.3	54.4 \pm 4.7	44.3 \pm 3.3	40.9 \pm 3.7	37.3 \pm 3.7
GPRGNN	LRBCD	76.2 \pm 7.9	73.1 \pm 7.9	66.1 \pm 11.5	60.5 \pm 11.7	56.4 \pm 12.4	53.4 \pm 12.9
	PRBCD	74.4 \pm 9.6	70.7 \pm 10.0	63.8 \pm 10.4	56.1 \pm 11.1	49.3 \pm 11.4	45.1 \pm 11.0
	EvA	71.9 \pm 10.4	65.0 \pm 11.4	50.1 \pm 11.4	41.8 \pm 14.0	37.2 \pm 14.6	35.2 \pm 15.7

to different subfields of machine learning. Each node is represented by a 1,433-dimensional bag-of-words (BoW) feature vector derived from the words in the titles and abstracts of the papers.

Citeseer. The CiteSeer dataset is also an academic citation network dataset consisting of 3,312 papers from 6 subfields of computer science and a total of 4,732 citation edges. Similar to Cora-ML, each paper as a node is represented by a BoW feature vector with a dimensionality of 3,703.

PubMed. The PubMed dataset is derived from a citation network of biomedical literature that contains 19,717 papers as nodes and 44,338 citation edges. Each paper is categorized into one of 3 classes based on its topic. The node features in PubMed are 500-dimensional vectors.

C.2 DETAILS OF MODELS

In the following sections, we detail the hyperparameters and architectural details for the models performed in this paper. The experimental configuration files, including all hyperparameters, will be made publicly available upon acceptance of the paper.

GCN. We utilize a two-layer GCN with 64 hidden units. A dropout rate of 0.5 is applied during training.

GAT. Our GAT model consists of two layers with 64 hidden units and a single attention head. During training, we apply a dropout rate of 0.5 to the hidden units, but no dropout is applied to the neighborhood.

APPNP. We use a two-layer MLP with 64 hidden units to encode the node attributes. We then apply generalized graph diffusion, using a transition matrix and coefficients $\gamma_K = (1 - \alpha)K$ and $\gamma_l = \alpha(1 - \alpha)l$ for $l < K$.

GPRGNN. Similar to APPNP, we employ a two-layer MLP with 64 hidden units for the predictive part. We use the symmetric normalized adjacency matrix with self-loops as the transition matrix and

Table 9: Classification accuracy (%) on the CoraML dataset under different attacks and adversarial training methods. The results are averaged over multiple runs with standard deviations.

Model	Adv. Tr.	Attack	ϵ						
			0.01	0.02	0.05	0.10	0.15	0.20	
GCN	None	LRBCD	78.51 \pm 1.56	75.94 \pm 1.54	71.10 \pm 1.16	64.41 \pm 1.65	60.14 \pm 1.73	57.37 \pm 1.45	
		PRBCD	76.44 \pm 1.64	73.17 \pm 1.39	66.48 \pm 2.13	58.51 \pm 1.77	52.67 \pm 2.09	47.19 \pm 2.02	
		EvA	74.80 \pm 1.50	68.97 \pm 1.58	52.95 \pm 1.91	41.99 \pm 2.06	37.65 \pm 2.74	35.37 \pm 2.38	
	LRBCD	LRBCD	79.64 \pm 1.77	77.51 \pm 2.41	73.10 \pm 1.54	68.19 \pm 1.11	64.84 \pm 1.92	62.35 \pm 3.00	
		PRBCD	78.79 \pm 1.88	75.87 \pm 1.41	69.75 \pm 1.81	62.35 \pm 2.70	56.80 \pm 3.04	54.23 \pm 4.71	
		EvA	76.80 \pm 1.29	71.10 \pm 1.64	56.30 \pm 1.66	48.40 \pm 2.91	43.06 \pm 2.43	40.85 \pm 2.70	
	PRBCD	LRBCD	80.71 \pm 1.16	77.86 \pm 0.81	73.81 \pm 0.54	69.40 \pm 0.91	66.48 \pm 1.08	63.77 \pm 1.73	
		PRBCD	78.93 \pm 1.27	76.30 \pm 1.27	70.25 \pm 1.74	64.06 \pm 1.83	59.50 \pm 2.84	56.58 \pm 4.53	
		EvA	76.80 \pm 0.77	71.53 \pm 1.65	57.44 \pm 2.13	49.11 \pm 3.05	44.70 \pm 2.96	41.92 \pm 3.32	
	EvA	LRBCD	80.85 \pm 1.36	78.58 \pm 0.99	74.66 \pm 1.11	69.89 \pm 0.93	66.98 \pm 0.92	65.05 \pm 1.08	
		PRBCD	79.79 \pm 1.80	76.51 \pm 1.31	71.25 \pm 1.54	64.34 \pm 1.97	60.43 \pm 1.32	58.22 \pm 2.26	
		EvA	77.22 \pm 1.87	71.96 \pm 2.38	57.94 \pm 3.08	50.04 \pm 3.29	44.91 \pm 3.44	42.63 \pm 2.33	
	GPRGNN	None	LRBCD	77.51 \pm 2.81	74.80 \pm 3.08	68.83 \pm 4.20	62.56 \pm 4.69	59.07 \pm 5.98	55.66 \pm 6.99
			PRBCD	74.95 \pm 3.08	71.67 \pm 2.76	64.84 \pm 4.18	57.94 \pm 4.55	53.24 \pm 5.20	48.68 \pm 6.52
			EvA	72.53 \pm 4.11	66.83 \pm 4.54	51.53 \pm 5.57	42.21 \pm 8.52	37.01 \pm 9.84	34.52 \pm 9.83
LRBCD		LRBCD	81.57 \pm 2.58	79.72 \pm 2.22	75.59 \pm 2.31	71.32 \pm 2.20	68.97 \pm 2.10	66.69 \pm 2.25	
		PRBCD	80.71 \pm 2.61	78.51 \pm 2.29	72.88 \pm 2.38	66.90 \pm 1.95	61.78 \pm 1.99	57.51 \pm 3.72	
		EvA	78.79 \pm 2.69	72.95 \pm 2.67	63.42 \pm 3.15	56.58 \pm 4.68	52.88 \pm 5.61	49.96 \pm 5.75	
PRBCD		LRBCD	80.43 \pm 2.01	78.01 \pm 1.91	73.74 \pm 1.66	69.96 \pm 2.14	67.19 \pm 2.51	64.84 \pm 3.20	
		PRBCD	80.21 \pm 2.43	77.30 \pm 2.63	71.53 \pm 2.67	65.12 \pm 3.21	60.07 \pm 4.10	55.37 \pm 3.85	
		EvA	78.79 \pm 2.45	73.10 \pm 2.54	62.85 \pm 4.93	56.94 \pm 6.64	53.74 \pm 7.65	51.60 \pm 8.10	
EvA		LRBCD	79.64 \pm 0.89	76.44 \pm 0.68	72.95 \pm 1.04	69.04 \pm 1.26	67.05 \pm 1.46	65.48 \pm 1.88	
		PRBCD	78.51 \pm 0.60	75.87 \pm 1.32	70.32 \pm 0.89	64.91 \pm 1.14	59.57 \pm 1.75	56.16 \pm 1.62	
		EvA	76.51 \pm 0.44	70.96 \pm 0.41	60.85 \pm 3.07	54.73 \pm 3.99	50.25 \pm 5.57	48.83 \pm 5.95	

Table 10: Dataset Statistics

Dataset	Nodes	Edges	Features	Classes
Cora-ML	2,810	7,981	1,433	7
Citeseer	3,312	4,732	3,703	6
PubMed	19,717	44,338	500	3

randomly initialize the diffusion coefficients. We consider a total of $K = 10$ diffusion steps, with α set to 0.1. During training, we apply a dropout rate of 0.2 to the MLP, while no dropout is applied to the adjacency matrix. Unlike the method in Chien et al. (2021), we always learn the diffusion coefficients with weight decay, which acts as a regularization mechanism to prevent the coefficients from growing indefinitely.

SoftMedian GDC. We follow the default configuration from Geisler et al. (2023), using a temperature of $T = 0.2$ or the SoftMedian aggregation, with 64 hidden dimensions and a dropout rate of 0.5. We fix the Personalized PageRank diffusion coefficient to $\alpha = 0.15$ and apply a top $k = 64$ sparsification. During the attacks, the model remains fully differentiable, except for the sparsification of the propagation matrix.

MLP. We design the MLP following the prediction module of GPRGNN and APPNP, incorporating two layers with 64 hidden units. During training, we apply a dropout rate of 0.2 to the hidden layer.

C.3 HYPERPARAMETER SETUP

In EvA we set the capacity of the computation to the same as the population, this means that all perturbations within a population are in one combined inference. However, in some cases where the

graph is large (e.g. PubMed), we reduce this number. Table 11 shows the hyper-parameter selection in almost all experiments. We only change the population number in some experiments, like certificate attacks, to reduce the computation. E.g., in the certificate attack, the population is reduced by a factor of 10.

C.4 ATTACK HYPERPARAMETERS

To assess the robustness of GNNs, we utilize the following attacks and hyperparameters. Based on Geisler et al. (2023), we also select the tanh-margin loss as the attack objective.

PRBCD. We closely adhere to the setup outlined by Geisler et al. (2023). A block size of 500,000 is used with 500 training epochs. Afterward, the model state from the best epoch is restored, followed by 100 additional epochs with a decaying learning rate and no block resampling. Additionally, the learning rate is scaled according to δ and the block size, as recommended by Geisler et al. (2023).

LRBCD. The same block size of 500,000 is used with 500 training epochs. The learning rate is scaled based on δ and the block size, following the same approach as PRBCD. The local budget is consistently set as 0.5.

EvA. We set the population size to 1024 in most cases. Our mutation rate is 0.01, and increasing this number breaks the balance between exploration and exploitation, leading to less effective attacks. We run each attack for 500 iterations in most cases. In cases like certificate attacks, which are time-consuming, we reduce this number to 100. The details are summarized in Table 11.

PGA. For the PGA, we adopt the same setting as in Zhu et al. (2023). We use GCN as the surrogate model and tanhMarginMCE-0.5 as the loss type. The attack is configured with 1 greedy step, a pre-selection ratio of 0.1, and a selection ratio of 0.6. Additionally, the influence ratio is set to 0.8, with the selection policy based on node degree and margin.

Table 11: Hyper-parameters for PRBCD, LRBCD, and EvA

Hyper-parameter	PRBCD	LRBCD	Hyper-parameter	EvA
Epochs	500	500	No. Steps	500
Fine-tune Epochs	100	0	Mutation Rate	0.01
Keep Heuristic	WeightOnly	WeightOnly	Tournament Size	2
Search Space Size	500,000	500,000	Population Size	1,024
Loss Type	tanhMargin	tanh-Margin	No. Crossovers	30
Early Stopping	N/A	False	Mutation Method	Adaptive

D DETAILS ON NOVEL OBJECTIVES

Smoothing-based certificate. We define a randomized model as a convolution of the original model and a smoothing scheme. The smoothing scheme $\xi : \mathcal{X} \mapsto \mathcal{X}$ is a randomized function mapping the given input to a random nearby point. For graph structure, we use the sparse smoothing certificate (Bojchevski et al., 2020), which certifies whether within \mathcal{B}_{r_a, r_d} the prediction of the smooth model remains the same. Here r_a is the maximum number of possible additions, and r_d is the maximum number of edge deletions. The smoothing function is defined by two Bernoulli parameters p_+ , and p_- ; i.e. for each entity of \mathcal{A} , if it is zero, it will be toggled with p_+ probability and otherwise with p_- .

The robustness certificate also accesses the model f as a black box and defines a smooth model as $\bar{f}_y(\mathbf{x}) = \mathbb{E}[\mathbb{I}[f(\xi(\mathbf{x})) = y]]$ - each random sample \mathbf{x}' is one vote for class $f(\mathbf{x}')$ and \bar{f}_y is the proportion of votes for class y . Let $p = \bar{f}_y(\mathbf{x})$; the certificate finds a lower bound probability $\underline{p} \geq \min_{\tilde{\mathbf{x}} \in \mathcal{B}(\mathbf{x})} \bar{f}_y(\tilde{\mathbf{x}})$ and acts as a decision function $\mathbb{I}[\underline{p} \geq 0.5]$. In other words, the certificate returns yes, in case it is guaranteed that the smooth model will not return any value lower than 0.5 for class y within $n\mathcal{B}(\mathbf{x})$. For further details about how to compute the certificate, see (Bojchevski et al., 2020).

Adaptive sampling for certificate attack. Statistical rigor is not a necessity while attacking the certificate. Therefore, while attacking, we can reduce the cost of resampling by only resampling the

subset of the graph that was perturbed. In other words, we initialize the search by computing samples $\mathbf{A}_1, \dots, \mathbf{A}_m$, and for each perturbation $\tilde{\mathbf{A}}$ we only resample the edges in $\mathbf{A} \Delta \tilde{\mathbf{A}}$. For each edge in that set, if the edge was added via the perturbations, we resample m Bernoulli variable with p_- , and otherwise p_+ . We substitute those samples in the same entry of $\mathbf{A}_1, \dots, \mathbf{A}_m$, and by running this process $|\delta|$ times, we assume that $\tilde{\mathbf{A}}_1, \dots, \tilde{\mathbf{A}}_m$ are representative as a new set of m samples for $\tilde{\mathbf{A}}$. This adaptive sampling reduces the number of random computations from $m \cdot n^2$ to $m \cdot |\delta|$, which is significantly lower. Surely, to evaluate the final perturbation (the reported effectiveness), we don't use this approach, as it is statistically flawed and only applicable to reduce the computation during the attack.

D.1 FITNESS FUNCTION

For a targeted attack, since the objective space is limited to the set $\{0, 1\}$, the sensitivity is very low. As a result, the method becomes equivalent to a random search. The zero-one fitness function means that all individuals with different perturbations receive the same score, causing the algorithm to behave more like a random search, as ties in each tournament are broken randomly. Secondly, since only one individual with a score of one is sufficient to halt the algorithm, all elite populations before success have scores of zero, which again results in random selection from them.

D.2 PERFORMANCE ON ARXIV

To demonstrate the scalability of our attack on larger datasets, we also present results for the Arxiv dataset. For this, we consider two realistic scenarios. In the first scenario, similar to the previous one, the attacker has access to modify a limited number of edges. We provide results for three values of epsilon: 0.1%, 0.5%, 1.0%. Table 12 summarizes the results for this scenario. EvA outperforms PRBCD for smaller budgets and achieves comparable performance with a 1% budget, which could be further improved by scaling and increasing computational resources.

Table 12: Comparison of PRBCD and EvA performance for varying ϵ values

Method	Clean	0.1%	0.5%	1%
PRBCD	70.53	69.83	68.64	66.27
EvA	70.53	69.21	67.59	66.86

Alternatively, in a more practical scenario, the attacker compromise a subset of nodes (e.g., 1,000 nodes) and get access to them, referred to as control nodes, and strategically target a specific group within the network (the target group). For example, in a social network, an attacker could purchase 1,000 user accounts and use them to influence the performance of other subgroups. For this experiment, we randomly sampled 1,000 nodes five times and also randomly selected 1,500 nodes 5 times as target group nodes. We then ran EvA and PRBCD and reported the average results in Table 13. Our method outperforms PRBCD in this scenario as well.

In summary, we demonstrate that our attack can be effectively applied and that it outperforms previous state-of-the-art methods.

Table 13: Comparison of PRBCD and EvA performance for varying ϵ values using control nodes

Method	Clean	1%	5%
PRBCD		64.89	54.7
EvA		59.3	53.92

D.3 TIME ANALYSIS

We also present an ablation study to compare the time analysis by evaluating the memory and wall clock time between EvA and the PRBCD method. In this experiment, we evaluate EvA with different numbers of steps, population sizes, and parallel evaluations, while PRBCD is run with varying

1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079

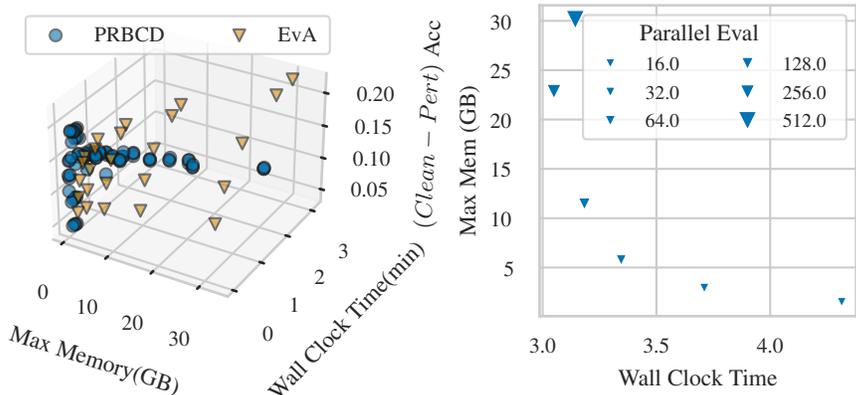


Figure 8: Comparing the memory usage between EvA and PRBCD

numbers of epochs and block sizes on the PubMed dataset. Fig. 8 (left) shows the results for EvA and PRBCD in terms of memory usage, wall clock time and method performance. Our method demonstrates comparable performance within the same level of wall clock time (less than a minute). Moreover, by increasing the wall clock time—through and memory either by a larger population size or more steps— EvA achieves additional benefits. Additionally, in Fig. 8 (right), we highlight how our framework provides a trade-off between time and memory for achieving the same level of accuracy by varying the number of parallel evaluations. For each point in the figure, we observe identical performance; however, the methods differ in memory usage due to different number of parallel evaluation, leading to variations in wall clock time.

D.4 COMPARISON WITH (DAI ET AL., 2018)

(Dai et al., 2018) proposed a practical black-box attack (PBA), dividing it into PBA-C (with access to logits - continuous) and PBA-D (access only to the labels - discrete). As stated in (Dai et al., 2018), a genetic algorithm for global attacks requires PBA-C because it relies on logits, with the fitness function being the negative log-likelihood. We demonstrate that EvA not only eliminates the need for logits but also performs even better by directly optimizing for accuracy rather than using log-likelihood. To compare our method with (Dai et al., 2018), we modified the algorithm’s fitness function and mutation mechanism to replicate the results reported in (Dai et al., 2018). This implementation retains scalability benefits, as it is also built upon our sparse encoded representation. Note here we re-implement Dai et al. (2018) in our sparse and parallelized framework. Their original implementation uses dense adjacency matrices and sequential evaluation and would achieve a significantly worse result within the same memory/run-time constraint. Even with our efficient re-implementation Dai et al. (2018) is significantly worse than ours. Table 14 provides the results for the CoraML dataset using the GCN architecture. EvA also significantly outperforms (Dai et al., 2018). Additionally, since our method is independent of gradients, we established the first attack on conformal prediction and certification. For conformal prediction, we attack coverage and set size where the latter criteria are not yet explored (to the best of our knowledge). Attacks tending to decrease certificate effectiveness are also under-explored in GNNs. In this work, we aim to achieve both attack on certified accuracy and certified ratio.

Attack Name	Clean	0.01	0.02	0.05	0.1	0.15	0.2
(Dai et al., 2018)	81.07 \pm 2.07	78.50 \pm 1.66	76.66 \pm 2.22	72.53 \pm 1.91	68.75 \pm 1.45	65.34 \pm 1.20	63.27 \pm 2.47
EvA	81.07 \pm 2.07	74.80 \pm 1.50	68.97 \pm 1.58	52.95 \pm 1.91	41.99 \pm 2.06	37.65 \pm 2.74	35.37 \pm 2.38

Table 14: Accuracy results of different attack methods under varying ϵ values.