
Friendly Noise against Adversarial Noise: A Powerful Defense against Data Poisoning Attacks

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Data poisoning attacks modify a subset of training examples by small adversarial
2 perturbations to change the prediction of certain test-time data. Existing defense
3 mechanisms are not desirable to deploy in practice, as they often drastically harm
4 the generalization performance, or are attack-specific and prohibitively slow to
5 apply. Here, we propose a simple but highly effective approach that unlike existing
6 methods breaks various types of poisoning attacks with the slightest drop in the
7 generalization performance. We make the key observation that attacks exploit sharp
8 loss regions to craft adversarial perturbations which can substantially alter exam-
9 ples' gradient or representations under small perturbations. To break poisoning
10 attacks, our approach comprises two components: an optimized friendly noise that
11 is generated to maximally perturb examples without degrading the performance,
12 and a random varying noise component. The first component takes examples farther
13 away from the sharp loss regions, and the second component smooths out the loss
14 landscape. The combination of both components builds a very light-weight but
15 extremely effective defense against the most powerful triggerless and backdoor
16 poisoning attacks, including Gradient Matching, Bulls-eye Polytope, and Sleeper
17 Agent. We show that our friendly noise is transferable to other architectures, and
18 adaptive attacks cannot break our defense due to its random noise component.

19 1 Introduction

20 Big datasets empower large over-parameterized deep learning systems. Such datasets are often
21 scraped from the internet or other public and user-provided sources. An adversary can easily insert a
22 subset of malicious examples into the data collected from public sources to harm the model's behavior.
23 As a result, deep learning systems trained on public data are extremely vulnerable to data poisoning
24 attacks, which modifies a subset of training examples under bounded adversarial perturbations with
25 the aim of changing the model's prediction on specific test-time examples. Powerful attacks generate
26 poisons that visually look innocent and are seemingly properly labeled [10, 15, 34]. This makes them
27 hard to detect even by expert observers. Hence, data poisoning attacks are arguably one of the most
28 concerning threats to deep learning systems [19].

29 Various types of poisoning attacks have been proposed to challenge and exploit the vulnerabilities of
30 deep learning systems. Backdoor data poisoning attacks add a fixed but not necessarily visible trigger
31 pattern to a subset of training data as well as the test-time target examples [13, 34, 39]. Triggerless
32 poisoning attacks add bounded perturbations to a subset of training examples to make them similar
33 to the adversarially labeled test-time target in the feature or gradient space [2, 12, 15, 32, 43]. In
34 both cases, training or fine tuning the model on the poisoned training data cause the model to
35 misclassify certain target examples at test-time.

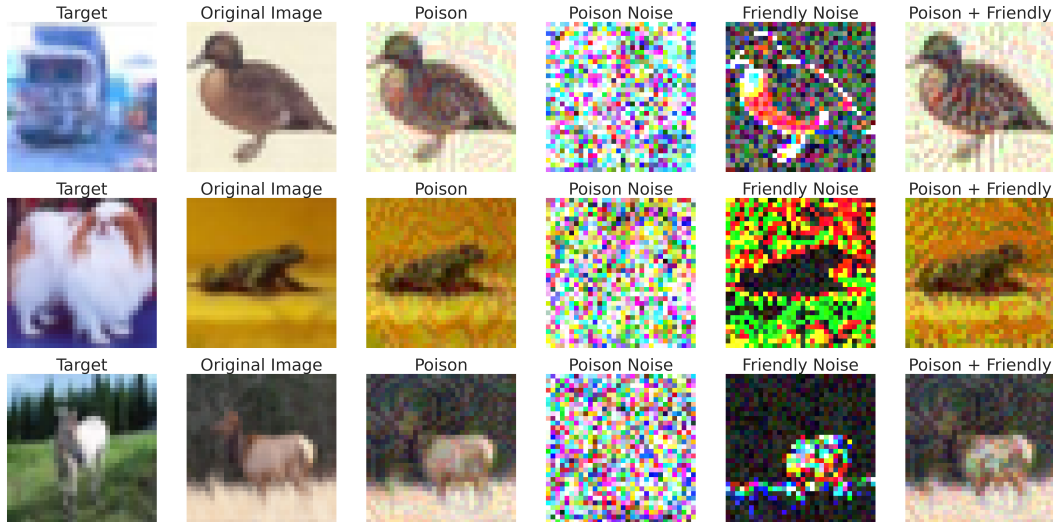


Figure 1: Qualitative Evaluation of Friendly Noise. Our optimized noise adds maximum allowed perturbation to the regions where network robustly learns and leaves other areas untouched (darker regions means less noise).

36 There has been sustained efforts to design effective defense mechanisms [1, 7, 11, 17, 25, 27, 33, 38].
 37 However, existing methods are highly impractical to be employed in real-world deep learning
 38 pipelines. Firstly, majority of the existing methods are attack specific and cannot protect the system
 39 against various types of data poisoning attacks [11, 27, 38]. Secondly, the provided protection is
 40 often in expense of significantly dropping the performance of the machine learning pipeline [1, 7, 25].
 41 Thirdly, existing methods are not effective in protecting the deep learning pipelines against adaptive
 42 attacks which can make more powerful poisons with the knowledge of the defense in place [17, 33].
 43 Finally, state-of-the-art defense methods are often so expensive that can hardly be applied to even
 44 medium-sized datasets [11, 27], and are ineffective in presence of larger number of poisons [7, 11, 27].

45 In this work, we propose a simple and powerful defense, namely Friendly Noise Defense (FRIENDS),
 46 against various types of data poisoning attacks. In particular, we make the following key observation:
 47 data poisoning attacks exploit sharp regions of the loss to craft adversarial perturbations which can
 48 substantially alter examples' gradient or representations under small perturbations. To effectively
 49 break poisoning attacks, our proposed method is composed of two noise components: first we find
 50 the maximum perturbation that can be added to every example without considerably changing the
 51 model's output. This fixed accuracy-friendly perturbation is found early in training and is transferable
 52 to other architectures. Then, we add a varying random noise in addition to the friendly perturbation to
 53 each example at every training iteration. Effectively, the friendly perturbation takes examples farther
 54 away from the sharp loss regions exploited by the attacker. The random noise component smooths
 55 out the loss landscape, and does not allow the crafted adversarial perturbation to match the target's
 56 gradient or representation closely. Despite being very light-weight, the combination of the two noisy
 57 components can effectively protect deep learning systems against various types of poisoning attacks,
 58 with minimum drop in the generalization performance.

59 We note that the random noise component of FRIENDS makes it extremely difficult for an adaptive
 60 attacker to break our defense. Adaptive attacks can bypass defenses by taking the defense mechanism
 61 into account when generating poisons. For FRIENDS, while an attacker may use the knowledge of
 62 the optimization procedure to bypass the friendly noise component, they need to take into account a
 63 prohibitively large number of random noise combinations when generating attacks. This makes it
 64 almost impossible for the attacker to ensure the effectiveness of an attack in presence of FRIENDS.

65 Through extensive experiments, we show that our light-weight method renders state-of-the-art
 66 targeted attacks, including Gradient Matching [10], Bullseye Polytope [2], Feature Collision [32],
 67 and Sleeper Agent [34] ineffective, with only a slight decrease in the performance. We also show
 68 that the optimized noise component generated based on a particular architecture can be applied to
 69 defend other architectures against data poisoning attacks. Therefore, it is easy to apply FRIENDS to
 70 real-world deep learning pipelines with minimal additional costs.

71 2 Related Work

72 **Targeted Data Poisoning Attacks.** Data poisoning attacks on deep networks have been explored
73 along two directions - triggered and triggerless attacks. Triggered attacks, or backdoor attacks, aim to
74 misclassify samples containing a ‘trigger’ patch as a pre-determined target class during inference time.
75 In the transfer learning or finetuning setting, earlier works of [8, 13, 23] relied on label modification or
76 unbounded image perturbations. These attacks are, however, easy to detect. Subsequently, [30, 34, 39]
77 introduced clean-label and visually imperceptible backdoor attacks. Recently, [34] proposed the first
78 clean-label hidden backdoor attack that is effective on victim models trained from scratch. Triggerless
79 data poisoning attacks aim to misclassify a given target as a pre-determined adversarial class, by
80 adding optimized bounded perturbations to a subset of training examples. Such attacks either optimize
81 for feature matching [2, 32, 43] or gradient matching [10] between poisoned and target images, or
82 use meta-learning to solve the poisoning problem directly via bilevel optimization [15].

83 **Defense Strategies.** Existing defenses against data poisoning can be divided into filtering and robust
84 training methods. Filtering methods detect outliers in feature space using thresholding [35] and
85 nearest neighbors [27], or activation space [7], or through decomposition of the feature covariance
86 matrix [38]. These defenses typically assume that only small subsets of the data are poisoned, hence
87 removing such points do not significantly harm generalization. In practice, this assumption may not
88 hold, and such defenses can be easily broken by increasing the number of poisons. Moreover, such
89 methods increase training time by orders of magnitudes, as the filtering step requires training the
90 model with poisons, followed by (usually expensive) filtering, and model retraining [7, 27, 35, 38].

91 Robust training methods apply randomized smoothing [41], strong data augmentation [4], or model
92 ensembling [21]. Other methods impose constraints on gradient magnitudes and directions [14],
93 detects and removes poisons with gradient ascent [22], or apply adversarial training [11, 25, 37].
94 Differentially private (DP) training methods have also been explored to defend against data poisoning
95 [1, 5, 16]. Robust training techniques usually involve a significant trade-off between generalization
96 and poison success rate [1, 14, 22, 25, 37], or are computationally very expensive [11, 25]. Compared
97 to augmentation-based and adversarial training methods, our method is simple, fast, and maintains
98 good generalization performance. Compared to data augmentation, the random noise component
99 of FRIENDS is considerably more effective in smoothing the loss landscape, due to its much larger
100 space of independent pixel-level transformations.

101 **Random and Adversarial Noise.** It is shown that small perturbations can result in large changes
102 in the output of a deep network [36]. Hence, application of random and adversarial noise has been
103 studied in various domains. In particular, [28] uses Gaussian noise to defend against query-based
104 black box attacks, and [29] shows that additive augmentations of Gaussian or Speckle noise is a
105 simple yet very strong baseline for robustness against image corruptions. Application of optimized
106 noise has been mainly studied in the context of adversarial training. [6] uses Generative Adversarial
107 Networks (GANs) to generate adversarial perturbations, and [24] relies on meta-learning to learn
108 a noise generator to defend against adversarial perturbations. Moreover, [26, 42] demonstrate the
109 transferability of adversarial perturbations across architectures and domains. In data poisoning, small
110 random noise generated from a particular distribution has been shown to be ineffective for breaking
111 attacks and harmful on the generalization performance [11]. In contrast, we show that random
112 noise combined with our proposed noise optimization approach, can make a highly effective defense
113 mechanism against data poisoning attacks and achieve a superior generalization performance.

114 3 FRIENDS: Friendly Noise Defense against Data Poisoning Attacks

115 Targeted data poisoning attacks modify a fraction of training data points by adding optimized
116 perturbations that are within an l_∞ -norm ξ -bound. The optimization is done with the objective of
117 changing the prediction of a target example x_t in the test set to an adversarial label y_{adv} . A small
118 perturbation bound ξ ensures that the poisoned examples remain visually similar to the original
119 (base) training data points. Poisons crafted by such attacks look innocent to human observer and are
120 seemingly labeled correctly. Hence, they are called clean-label attacks. Targeted clean label data
121 poisoning attacks can be formulated as the following bilevel optimization problem:

$$\min_{\delta \in \mathcal{C}} \mathcal{L}(x_t, y_{adv}, \theta(\delta)) \quad s.t. \quad \theta(\delta) = \arg \min_{\theta} \sum_{i \in V} \mathcal{L}(x_i + \delta_i, y_i, \theta), \quad (1)$$

Algorithm 1 Generating Friendly Noise

Require: Train dataset X , Model f_θ , LR η_{opt} , λ , small number T

```
for  $i \in [T]$  do  
   $\theta^i = \theta^{i-1} - \eta \nabla_\theta L(\theta^{i-1}, X)$  ▷ Train the model for a few epochs  
end for  
for  $x_i \in X$  do  
  Initialize noise  $\epsilon_i^0$  uniformly sampled from  $[-\epsilon_{init}, \epsilon_{init}]$   
  for  $t = 1$  to  $T$  do  
     $\epsilon_i^t = \epsilon_i^{t-1} - \eta_{opt} \nabla_\epsilon (D_{KL}(f_\theta(x_i + \epsilon_i^{t-1}) || f_\theta(x_i)) - \lambda \|\epsilon_i^{t-1}\|_2)$   
  end for  
  Store noise  $\epsilon_i = \epsilon_i^T$  for example  $x_i$   
end for
```

122 where $\mathcal{C} = \{\delta \in \mathbb{R}^{n \times m} : \|\delta\|_\infty \leq \xi, \delta_i = 0 \forall i \notin V_p\}$ is the constraint set defining the set of valid
123 poisons, V is the training data, and V_p is the set of poisoned training examples. To address the
124 above optimization problem, powerful poisoning attacks such as Meta Poison (MP) [15], Gradient
125 Matching(GM) [10], Bull-eyes Polytope (BP) [2], and Sleeper Agent [34] craft the poisons to mimic
126 the gradient (equivalently representation in transfer learning) of the adversarially labeled target, i.e.,

$$\nabla \mathcal{L}(x_t, y_{adv}, \theta) \approx \frac{1}{|V_p|} \sum_{i \in V_p} \nabla \mathcal{L}(x_i + \delta_i, y_i, \theta), \quad (2)$$

127 Minimizing the training loss on RHS of Eq.(1) also minimizes the adversarial loss on LHS of Eq. (1).

128 3.1 Powerful Poisons Fall in Sharp Loss Regions

129 Based on Eq. (2), we make the following key observation. To substantially change the gradient of a
130 training example $\nabla \mathcal{L}(x_t, y_{adv}, \theta) \approx \nabla \mathcal{L}(x_i + \delta_i, y_i, \theta)$, under bounded perturbation $\|x_i + \delta_i\|_\infty \leq$
131 $\|x_i\|_\infty + \xi$, the attacker needs to exploit the highly non-convex nature of the loss by finding sharp
132 regions in a ball of radius ξ around example x_i . If such regions can be found, the example can be
133 slightly modified to fall in the sharp region and further optimized there to match the target gradient.
134 Fig. 2(c) shows that the undefended model’s output around every poison has a large variance. Such
135 non-linearities can be exploited by attackers to craft poisons effectively.

136 As poisons rely on sharp loss regions to match a particular gradient or representation, they are highly
137 sensitive to small perturbations. Indeed, slightly perturbing the poisons considerably change their gra-
138 dient and make them ineffective. The main idea behind our friendly noise defense method, FRIENDS,
139 is to maximally perturb the training examples to make the poisons ineffective. However, perturbing
140 training examples should be done in a way that does not harm the generalization performance of
141 the model. To address this, our method is composed of two components: First, we find maximum
142 perturbation that can be added to every training example without changing its prediction. Intuitively,
143 this pulls examples far away from sharp loss regions that may have been exploited by the attacker. To
144 further break the attack, we add varying random noise to every example during the training. This
145 smooths out the loss and consecutively changes the poisons gradient and makes the crafted poisons
146 ineffective. Below, we discuss each component in more details.

147 3.2 Optimizing the Friendly Noise: Taking Poisons away from Sharp Regions

148 The first component of our method finds the maximum perturbation that can be added to every
149 example without considerably changing the model’s output. To do so, for every example x_i we
150 optimize for the largest noise ϵ_i within an l_∞ norm ζ -bound that results in a similar class probabilities,
151 measured by KL-divergence. Formally for each example x_i , we find perturbation ϵ_i as follows:

$$\epsilon_i = \arg \min_{\epsilon: \|\epsilon\|_\infty \leq \zeta} D_{KL}(f_\theta(x_i + \epsilon) || f_\theta(x_i)) - \lambda \|\epsilon\|_2 \quad (3)$$

152 We generate a fixed accuracy-friendly perturbation for every data point by solving problem (3)
153 once, using a few Stochastic Gradient Descent (SGD) steps. There is a trade-off between the
154 time of generating the optimized perturbations and their effectiveness. In particular, the optimized
155 perturbations need to be generated and added to the training examples early in training, before the

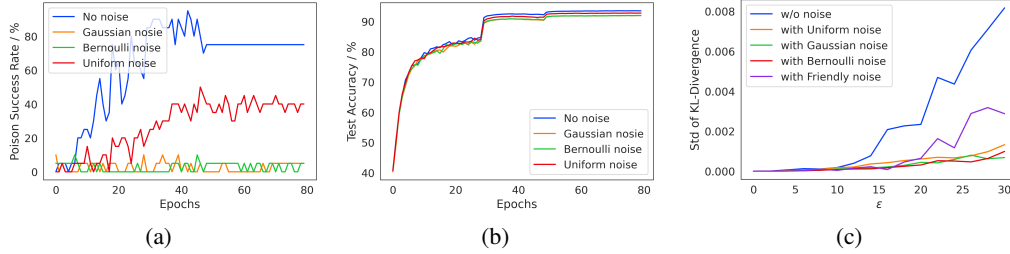


Figure 2: (a) Number of poisoned datasets vs epochs. It takes at least a few epochs for the poisons to have an effect. (b) On the other hand, generalization of the predictions and features learnt increases over time, as measured by error on the test set. (c) Standard deviation of KL-divergence (y-axis) between predictions of training examples before and after adding friendly perturbations in Eq. (3) and random noise sampled from various distributions. Standard deviation is calculated over 10 randomly sampled points in an ϵ balls (x-axis) around every training example.

156 attack succeeds (Fig. 2(a)). At the same time for the perturbations to be effective, they should be
 157 generated after the decision boundary is shaped (Fig. 2(b)). We found that training for as few as
 158 5 epochs before solving the optimization problem (3) yields effective perturbations that reduce the
 159 attack success rate without harming the model’s performance. The pseudocode can be found in Alg.1.

160 To better understand the effect of our friendly noise, we first look at the histogram of the noise added to
 161 every pixel in the training data. Fig. 3(a) shows that our method mainly targets certain pixels in every
 162 image by adding the maximum amount of perturbation, and leaves the rest of the image untouched.
 163 Certain semantic regions have been shown to be much more robust to perturbations [3]. Perturbing
 164 the more robust areas does not considerably change the model’s behavior and training dynamics. On
 165 the other hand, the least robust areas are very sensitive to perturbations, hence small amounts of noise
 166 in such areas result in a relatively large change in the model’s behaviour. We visualize the poisons
 167 before and after adding our friendly noise in Fig. 1. We see that our friendly noise successfully targets
 168 certain areas in every image that are robustly learned, by adding the maximum possible perturbation.

169 Intuitively, perturbing the more robust areas pulls examples far away from the sharp regions of the
 170 loss, and perturbing the less robust regions move them closer to the sharp regions. Our friendly
 171 perturbation successfully pulls the poisons far away from the sharp regions that may have been
 172 utilized by the attack. By moving along the level curve of the loss, the important features used for
 173 classifications are preserved and hence model predictions remain unchanged, but the adversarial
 174 poison perturbations can no longer closely match the target gradient or representation. Hence, the
 175 crafted adversarial perturbations added to the examples augmented with our friendly noise cannot
 176 poison the model. Hence, our method reduces the attack success rate while preserves the training
 177 dynamics and ensures minimum drop in the test accuracy.

178 Next, we discuss how adding a random variable random noise can further improve the model’s
 179 robustness against poisoning attacks.

180 3.3 Adding Random Noise: Smoothing the Loss

181 As discussed, our friendly perturbation mainly targets the robust areas of the image and pulls examples
 182 far away from the sharp loss regions, without considerably changing the training dynamics. To further
 183 improve model’s robustness against poisoning attacks, we add a variable random noise to the training
 184 examples at every training iteration. Effectively, adding the variable random noise smooths out the
 185 loss landscape. In doing so, the optimized poison perturbations cannot match the target gradient, as
 186 on a smooth manifold the gradients do not change considerably in a small ball around every example.
 187 Consequently, adding the random noise considerably drops the attack success rate. The pseudocode
 188 of our defense, FRIENDS, can be found in Alg. 2.

189 The small varying random noise can be sampled from various distributions. In particular, we sample
 190 from 3 different random noise distributions: (1) Bernoulli: noise is randomly sampled from $\{-\mu, \mu\}$,
 191 (2) Uniform: noise is randomly sampled from $[-\mu, \mu]$, and (3) Gaussian: noise is sampled from
 192 the normal distribution $\mathcal{N}(0, \mu)$. Fig. 3 compares the distribution of random noise sampled from

Algorithm 2 Training with FRIENDS

Require: Train dataset X , Random Noise Distribution A , Epoch to start defense def_epoch

```
Run Algorithm 1 to generate  $\{\epsilon_i\}_{i=1}^{|X|}$ 
for  $i = def\_epoch$  to  $n\_epochs$  do
  for  $x_i \in X$  do
    Sample random noise  $\mu_i \sim A$ 
    Set  $\hat{x}_i = x_i + \epsilon_i + \mu_i$  and add to dataset  $\hat{X}$ 
  end for
   $\theta^i = \theta^{i-1} - \eta \nabla_{\theta} L(\theta^{i-1}, \hat{X})$  ▷ SGD update step with new dataset  $\hat{X}$ 
end for
```

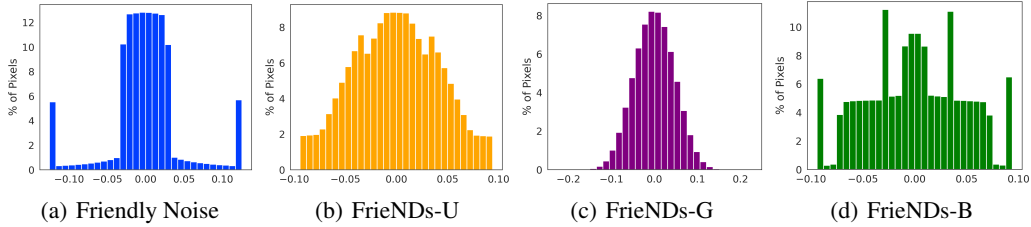


Figure 3: Histogram of our method with different types of random noises For (a), we set $\zeta = 32$. For (b)-(d), we set $\zeta = 16, \mu = 16$.

193 these distributions combined with the optimized perturbation obtained from Eq. (3). We can see that
194 uniform noise perturbs all the pixels similarly, and hence small amount of uniform noise does not
195 harm the model’s performance but cannot effectively breaks poisons. Larger uniform noise, however,
196 has a larger effect on the model’s performance. Bernoulli and Gaussian noise target a smaller number
197 of pixels but add a larger perturbation to them. Hence, they are more effective in reducing the attack
198 success rate, but they harm the test accuracy more as they the larger perturbation may be added to the
199 more sensitive areas. Figures 3(b) to 3(d) shows the distribution of random noise combined with our
200 optimized noise added to different pixels. We observe that random noise is added to regions where
201 friendly noise is less dominant, hence resulting in a significant perturbation to an overall greater
202 number of pixels to negate poisoning attacks.

203 Fig. 2(c) shows that the standard deviation of the model output in a ball around every example
204 becomes considerably smaller after adding each component of our defense. This clearly demonstrates
205 the effect of our defense and explains its effectiveness.

206 3.4 Adaptive attacks

207 Adaptive attacks can respond to a novel defense algorithm when the attacker is aware of the defense.
208 If the defense algorithm is known to the attacker beforehand, the attacker can generate more powerful
209 poisons by taking into account the specific defense in place. For example, Gradient Matching [10]
210 and Sleeper Agent [34] demonstrated that by including augmented examples as well as original
211 examples during poison generation in Eq. (2), they can obtain robustness against standard data
212 augmentations like crops and flips, when the augmentation technique is preempted by the attacker.
213 For FRIENDS, the prohibitively large search space of random noise permutations and its pixel-wise
214 independence property makes it nearly impossible for adaptive attacks to break. That is, while
215 an attacker may use the knowledge of the time and optimization procedure to bypass the friendly
216 perturbation component of FRIENDS, they needs to take into account a prohibitively large number
217 of random noise combinations to bypass the random noise component. For example, for a fixed
218 Bernoulli noise, with p pixels there are 2^p combinations that an attacker should take into account to
219 ensure the poisons’ effectiveness. For real images this becomes prohibitively expensive. Similarly for
220 a fixed Gaussian and uniform noise, there are an infinite number of combinations to be considered
221 during the poison generation to ensure attack’s robustness. Note that our method applies a varying
222 random noise at every iteration, which also needs to be taken into account by the attacker. This makes
223 FRIENDS robust against adaptive attacks.

Table 1: Baselines - Against Gradient Matching eps=16, 80 epochs. For trials with all equal outcomes, we report worst-case error estimate 5.59%

Defense	Poison Acc	Test Acc	Time (HH:MM)
AP-0.25 [11]	15.00% ($\pm 2.85\%$)	93.27% ($\pm 0.00\%$)	02:39
AP-0.5 [11]	10.00% ($\pm 2.01\%$)	92.83% ($\pm 0.00\%$)	03:41
AP-0.75 [11]	0.00% ($\pm 5.59\%$)	91.29% ($\pm 0.00\%$)	04:30
DeepKNN [27]	75.00% ($\pm 4.19\%$)	93.72% ($\pm 0.24\%$)	02:55
Adversarial Training [25]	60.00% ($\pm 5.37\%$)	92.03% ($\pm 0.31\%$)	02:37
Activation Clustering [7]	45.00% ($\pm 5.53\%$)	87.69% ($\pm 0.50\%$)	01:01
Diff. Priv. SGD [14]	5.00% ($\pm 1.06\%$)	75.70% ($\pm 1.19\%$)	00:38
Friendly Noise	10.00% ($\pm 2.01\%$)	91.73% ($\pm 0.30\%$)	00:37
FRIENDS-U	5.00% ($\pm 2.01\%$)	91.91% ($\pm 0.28\%$)	00:37
FRIENDS-B	0.00% ($\pm 5.59\%$)	91.52% ($\pm 0.28\%$)	00:37
FRIENDS-G	0.00% ($\pm 5.59\%$)	91.50% ($\pm 0.25\%$)	00:37

Table 2: Comparisons with Sleeper Agent defenses averaged over 24 datasets (40 epoch setting)

Defense	Poison Acc	Test Acc
None	83.48% ($\pm 7.58\%$)	91.56% ($\pm 0.19\%$)
Spectral Signatures [38]	37.17% ($\pm 10.10\%$)	89.94% ($\pm 0.19\%$)
Activation Clustering [7]	15.17% ($\pm 5.38\%$)	72.38% ($\pm 0.48\%$)
Diff. Priv. SGD [14]	13.14% ($\pm 4.49\%$)	70.00% ($\pm 0.17\%$)
Strong Augmentation [5]	69.75% ($\pm 10.77\%$)	91.32% ($\pm 0.12\%$)
STRIP [9]	62.68% ($\pm 4.90\%$)	92.23% ($\pm 0.05\%$)
NeuralCleanse [40]	85.11% ($\pm 5.04\%$)	92.26% ($\pm 0.06\%$)
FRIENDS-B	21.52% ($\pm 8.39\%$)	89.76% ($\pm 0.30\%$)
FRIENDS-G	22.99% ($\pm 8.59\%$)	89.87% ($\pm 0.31\%$)
FRIENDS-U	34.53% ($\pm 9.71\%$)	90.36% ($\pm 0.38\%$)

224 4 Experiments

225 4.1 Implementation details

226 We evaluate our defense method against both triggerless data poisoning and backdoor attacks under
 227 two attack settings - training from scratch and transfer learning. Following the works of [10, 11, 31],
 228 we evaluate our method primarily on CIFAR-10, ResNet-18. We also normalize and augment training
 229 images with default CIFAR-10 augmentations as used in [11]. For all models trained from scratch, we
 230 use a learning rate starting at 0.1 and decaying by a factor of 10 at epochs 30, 50, and 70. For transfer
 231 learning, we decay learning rate at epochs 15, 25, and 35. When applying our method, we clamp the
 232 generated friendly perturbations using $\zeta = 16$, and add bounded random noise. For the random noise
 233 component, we set $\mu = 16$ in our experiments. We also normalize the image as a pre-processing
 234 step. We optimize friendly perturbations using SGD with momentum 0.9 and Nesterov acceleration,
 235 perform a hyperparameter search along $LR = \{10, 20, 50, 100\}$ and $\lambda = \{1, 10\}$, and optimize each
 236 batch of 128 samples for 20 epochs. Following previous works, we report poison success rate (or
 237 poison accuracy) as the percentage of datasets poisoned at the end of training. We run all experiments
 238 and timings on an NVIDIA A40 GPU.

239 4.2 From-Scratch Setting

240 First, we evaluate our method on poisoning attacks targeted towards victim models trained from
 241 scratch. Such attack assumes a gray-box scenario, where attackers have knowledge of the victim
 242 architecture, but have no knowledge of the specific initialization of the victim’s model. Similar to
 243 the settings used in [31], which proposes a standardized benchmark for backdoor and data poisoning
 244 attacks, benchmark settings, we generate poisoning attacks by selecting 1% of training examples as
 245 poisons, which are perturbed within the l_∞ ball of some radius ξ . Unless otherwise specified, we set
 246 $\xi = 16$. The victim model is initialized with the same architecture targeted by the attack based on a

Table 3: Against different data poisoning attacks. Here, we use FRIENDS-B as the defense method. Note: Baseline metapoison is ran without default augmentations, following settings used in [11, 15]

Attack	Scenario	Undefended		Defended	
		Posion Acc	Test Acc	Poison Acc	Test Acc
Gradient Matching ($\xi = 8$)	From-scratch	50.00%	93.55%	0.00%	91.55%
Gradient Matching ($\xi = 16$)	From-scratch	75.00%	93.50%	0.00%	91.52%
Metapoison ($\xi = 8$)	From-scratch	45.00%	87.61%	20.00%	90.82%
Bullseye Polytope	Transfer	100.00%	92.13%	35.00%	79.35%
Poison Frogs	Transfer	100.00%	92.12%	30.00%	79.07%
Sleeper Agent	From-scratch	91.72%	93.36%	31.20%	91.31%

Table 4: Ablation study on random noise components of FRIENDS using Gradient Matching attack ($\xi = 16$). We set $\zeta = 32$ for Friendly Noise, and $\mu = 32$ for Noise Only. For experiments on FRIENDS, we set $\zeta = 16$, $\mu = 16$ to combine each component proportionately.

No Def.		Friendly Noise		Noise Type	Noise Only		FRIENDS	
P. Acc.	Test Acc.	P. Acc.	Test Acc.		P. Acc.	Test Acc.	P. Acc.	Test Acc.
75.00	93.50	10.00	91.73	Gaussian	0.00	89.46	0.00	91.50
				Bernoulli	5.00	89.31	0.00	91.52
				Uniform	0.00	91.61	5.00	91.91

247 different random seed, and is trained on the poisoned dataset using SGD. When applying FRIENDS,
 248 we set $def_epoch = 5$, and train only with random noise for the first 5 epochs.

249 4.2.1 Baseline Comparison and Ablation Study

250 We evaluate our method and baseline defenses against the Witches’ Brew, or Gradient Matching,
 251 attack [10]. It is the current state-of-the-art among data poisoning attacks when applied to the
 252 from-scratch setting, and is adapted to be effective against data augmentation and differential privacy
 253 [11]. We follow the settings proposed by [11], under which we generate 20 different attack datasets
 254 for ResNet-18 trained on CIFAR10 with a 1% budget bounded by $\xi = 16$, with a slight modification -
 255 while [11] uses 40 epochs for training, we use 80 epochs to show that our method easily scales to
 256 real-world training pipelines. This is because 40 epochs of training only yields 92.01% test error,
 257 while 80 epochs yield 93.50%. In Tab. 1, we show that we outperform state-of-the-art defenses
 258 [7, 11, 14, 25, 27]. We achieve the same 0.00% poison success rate with 91.52% test accuracy, an
 259 improvement over of 0.23% over state-of-the-art [11] which yields 91.29% test accuracy at the same
 260 poison success rate. Most importantly, FRIENDS completes in 37 mins, 7.3x faster than [11] which
 261 completes in 4.5hrs. We also strongly outperform other baseline defense methods simultaneously in
 262 all three metrics - poison success rate, test accuracy, and runtime. In Tab. 3, we show that FRIENDS
 263 also effectively defends against MetaPoison [15], reducing the poison success rate from 45.00% to
 264 20.00% with an accuracy gain from 87.61% to 90.82% resulted from applying augmentations.

265 We further show that our approach is effective against backdoor attacks, in particular, against the
 266 Sleeper Agent attack [34]. Sleeper Agent is the current state-of-the-art clean-label backdoor attack,
 267 and the only such attack shown to be effective in from-scratch settings. Following their evaluation
 268 protocol, we generate 24 poisoned datasets with $\xi = 16$, and evaluate our defense by training 24
 269 victim models respectively for 40 epochs and test poison success rate on 1000 target backdoor images
 270 per dataset. We compare our method against other defenses evaluated by [34] in Tab. 2. Here,
 271 FRIENDS successfully defends against [34] by reducing poison accuracy from 83.48% to 21.52%
 272 with only a small drop in test accuracy from 91.56% to 89.76%. We outperform the next best method,
 273 Spectral Signatures [38], by lowering poison accuracy by 14.18% while maintaining similar test
 274 accuracy. [14] achieves the lowest poison success rate at 13.14%, but causes a significant drop in test
 275 accuracy to 70.00%, and [9] achieves 92.26% test accuracy but suffers from 62.68% poison accuracy.

276 We also perform an ablation on each components of FRIENDS in Tab. 4. We show that naively
 277 applying Friendly Noise ($\zeta = 32$) yields a high poison success rate of 10%. On the other hand,
 278 applying random noise ($\mu = 32$) yields low poison success rates but also results in a significant
 279 test accuracy tradeoff (e.g. $> 4.0\%$ drop for Gaussian and Bernoulli noise). Here, we show that

Table 5: Transferability between different architectures

Method	Poison Acc	Test Acc
FRIENDS-B (ResNet18)	0%	91.52%
FRIENDS-B (AlexNet -> ResNet18)	0%	91.27%
FRIENDS-B (LeNet -> ResNet18)	0%	91.39%

280 applying FRIENDS by proportionately combining friendly noise ($\zeta = 16$) with each of the random
 281 noise components ($\mu = 16$) maintains high test accuracy (i.e. only 2.0% drop) while keeping poison
 282 success rate close to 0.

283 4.3 Transfer learning

284 Next, we evaluate our method on data poisoning and backdoor attacks designed for the transfer
 285 learning scenario [2, 32]. Here, the attacks are crafted based on a pretrained network with the goal of
 286 achieving poisoning when transfer learning is performed using the generated poisoned dataset. For
 287 the transfer learning scenario used in poisoning benchmarks [11, 31], the linear layer (classifier) of
 288 the pretrained model is re-initialized and trained with the poisoned dataset, while other layers (feature
 289 extractor) remain fixed during the training. Similar to the from-scratch setting, attacks are limited to
 290 a budget of 1% and $\xi = 16$. However, we generate FRIENDS at the beginning of training instead
 291 of after 5 training epochs, since the feature extractor is already initialized. We note that this is not
 292 the true transfer learning setting, since the pretraining and transfer learning datasets are the same.
 293 However, as [10, 11] note, this presents an effective worst-case scenario to evaluate poisoning attacks.
 294 We show that in Tab. 3 that even in such cases, we reduce poison success rate from 100% to 35% for
 295 the Bullseye Polytope attack [2], and from 100% to 30% for the Poison Frogs attack [32].

296 4.4 Transferability across Architectures

297 We show that perturbations generated by FRIENDS are transferable across architectures. In Tab. 5,
 298 we show using Gradient Matching $\xi = 16$ that FRIENDS optimized using smaller architectures, in
 299 particular AlexNet [18] and LeNet [20], can be directly used for larger architectures like ResNet18.
 300 This presents a significant advantage in terms of computational costs, since FRIENDS can be
 301 generated using smaller, and hence faster, models. Crucially, this makes the generated friendly noise
 302 free to be directly applied to (much larger) architectures.

303 5 Conclusion

304 We proposed a simple and highly effective defense mechanism, FRIENDS, that protects deep learning
 305 pipelines against various types of poisoning attacks. Our defense is built on the key observation
 306 that poisoning attacks exploit sharp regions of the loss to craft adversarial perturbations that when
 307 added to an example, can substantially alter its gradient or representations. FRIENDS relies on two
 308 components to break the poisons: an accuracy-friendly perturbation that is generated to maximally
 309 perturb examples without degrading the performance, and a random varying noise component. The
 310 first component takes examples farther away from the sharp loss regions, and the second component
 311 smooths out the loss landscape. Both components combined together build a very light-weight
 312 but highly effective defense against the most powerful triggerless and backdoor poisoning attacks,
 313 including Gradient Matching, Bull-eyes Polytope, Poison Frogs, and Sleeper Agent, in transfer
 314 learning or training from scratch scenarios. FRIENDS is impossible to break with adaptive attacks
 315 and our friendly noise can be transferred to other architecture. This makes it almost free to apply to
 316 real-world deep learning pipelines. Our defense is particularly targeted towards clean-label poisoning
 317 attacks that are generated under bounded perturbations. Such settings are the most difficult to defend,
 318 as generated poisons can easily fool even an expert observer. In contrast, unbounded attacks can be
 319 easily detected by manual or automated filtering mechanisms, through a single pass over the dataset.

320 References

321 [1] Martin Abadi, Andy Chu, Ian Goodfellow, H Brendan McMahan, Ilya Mironov, Kunal Talwar,
 322 and Li Zhang. Deep learning with differential privacy. In *Proceedings of the 2016 ACM SIGSAC*

- 323 *conference on computer and communications security*, pages 308–318, 2016.
- 324 [2] Hojjat Aghakhani, Dongyu Meng, Yu-Xiang Wang, Christopher Kruegel, and Giovanni Vigna.
325 Bullseye polytope: A scalable clean-label poisoning attack with improved transferability. In
326 *2021 IEEE European Symposium on Security and Privacy (EuroS&P)*, pages 159–178. IEEE,
327 2021.
- 328 [3] Zachary Berger, Parth Agrawal, Tian Yu Liu, Stefano Soatto, and Alex Wong. Stereoscopic uni-
329 versal perturbations across different architectures and datasets. *arXiv preprint arXiv:2112.06116*,
330 2021.
- 331 [4] Eitan Borgnia, Valeriia Cherepanova, Liam Fowl, Amin Ghiasi, Jonas Geiping, Micah Goldblum,
332 Tom Goldstein, and Arjun Gupta. Strong data augmentation sanitizes poisoning and backdoor
333 attacks without an accuracy tradeoff. In *ICASSP 2021-2021 IEEE International Conference on*
334 *Acoustics, Speech and Signal Processing (ICASSP)*, pages 3855–3859. IEEE, 2021.
- 335 [5] Eitan Borgnia, Jonas Geiping, Valeriia Cherepanova, Liam Fowl, Arjun Gupta, Amin Ghiasi,
336 Furong Huang, Micah Goldblum, and Tom Goldstein. Dp-instahide: Provably defusing poi-
337 soning and backdoor attacks with differentially private data augmentations. *arXiv preprint*
338 *arXiv:2103.02079*, 2021.
- 339 [6] Yuanhao Cai, Xiaowan Hu, Haoqian Wang, Yulun Zhang, Hanspeter Pfister, and Donglai Wei.
340 Learning to generate realistic noisy images via pixel-level noise-aware adversarial training.
341 *Advances in Neural Information Processing Systems*, 34, 2021.
- 342 [7] Bryant Chen, Wilka Carvalho, Nathalie Baracaldo, Heiko Ludwig, Benjamin Edwards, Taesung
343 Lee, Ian Molloy, and Biplav Srivastava. Detecting backdoor attacks on deep neural networks by
344 activation clustering. In *SafeAI@ AAAI*, 2019.
- 345 [8] Xinyun Chen, Chang Liu, Bo Li, Kimberly Lu, and Dawn Song. Targeted backdoor attacks on
346 deep learning systems using data poisoning. *arXiv preprint arXiv:1712.05526*, 2017.
- 347 [9] Yansong Gao, Change Xu, Derui Wang, Shiping Chen, Damith C. Ranasinghe, and Surya Nepal.
348 Strip: A defence against trojan attacks on deep neural networks. In *Proceedings of the 35th*
349 *Annual Computer Security Applications Conference, ACSAC '19*, page 113–125, New York,
350 NY, USA, 2019. Association for Computing Machinery.
- 351 [10] Jonas Geiping, Liam Fowl, W Ronny Huang, Wojciech Czaja, Gavin Taylor, Michael Moeller,
352 and Tom Goldstein. Witches’ brew: Industrial scale data poisoning via gradient matching. *arXiv*
353 *preprint arXiv:2009.02276*, 2020.
- 354 [11] Jonas Geiping, Liam Fowl, Gowthami Somepalli, Micah Goldblum, Michael Moeller, and Tom
355 Goldstein. What doesn’t kill you makes you robust (er): Adversarial training against poisons
356 and backdoors. 2021.
- 357 [12] Jonas Geiping, Liam H Fowl, W. Ronny Huang, Wojciech Czaja, Gavin Taylor, Michael Moeller,
358 and Tom Goldstein. Witches’ brew: Industrial scale data poisoning via gradient matching. In
359 *International Conference on Learning Representations*, 2021.
- 360 [13] Tianyu Gu, Brendan Dolan-Gavitt, and Siddharth Garg. Badnets: Identifying vulnerabilities in
361 the machine learning model supply chain. *arXiv preprint arXiv:1708.06733*, 2017.
- 362 [14] Sanghyun Hong, Varun Chandrasekaran, Yiğitcan Kaya, Tudor Dumitraş, and Nicolas Papernot.
363 On the effectiveness of mitigating data poisoning attacks with gradient shaping. *arXiv preprint*
364 *arXiv:2002.11497*, 2020.
- 365 [15] W Ronny Huang, Jonas Geiping, Liam Fowl, Gavin Taylor, and Tom Goldstein. Metapoi-
366 son: Practical general-purpose clean-label data poisoning. *Advances in Neural Information*
367 *Processing Systems*, 33, 2020.
- 368 [16] Bargav Jayaraman and David Evans. Evaluating differentially private machine learning in
369 practice. In *28th {USENIX} Security Symposium ({USENIX} Security 19)*, pages 1895–1912,
370 2019.

- 371 [17] Pang Wei Koh, Jacob Steinhardt, and Percy Liang. Stronger data poisoning attacks break data
372 sanitization defenses. *arXiv preprint arXiv:1811.00741*, 2018.
- 373 [18] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep
374 convolutional neural networks. In F. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger,
375 editors, *Advances in Neural Information Processing Systems 25*, pages 1097–1105. Curran
376 Associates, Inc., 2012.
- 377 [19] Ram Shankar Siva Kumar, Magnus Nyström, John Lambert, Andrew Marshall, Mario Goertzel,
378 Andi Comissioneru, Matt Swann, and Sharon Xia. Adversarial machine learning – industry
379 perspectives, 2020.
- 380 [20] Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document
381 recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- 382 [21] Alexander Levine and Soheil Feizi. Deep partition aggregation: Provable defenses against
383 general poisoning attacks. In *International Conference on Learning Representations*, 2020.
- 384 [22] Yige Li, Xixiang Lyu, Nodens Koren, Lingjuan Lyu, Bo Li, and Xingjun Ma. Anti-backdoor
385 learning: Training clean models on poisoned data. *Advances in Neural Information Processing
386 Systems*, 34, 2021.
- 387 [23] Yingqi Liu, Shiqing Ma, Yousra Aafer, Wen-Chuan Lee, Juan Zhai, Weihang Wang, and
388 Xiangyu Zhang. Trojaning attack on neural networks. 2017.
- 389 [24] Divyam Madaan, Jinwoo Shin, and Sung Ju Hwang. Learning to generate noise for multi-attack
390 robustness. In *International Conference on Machine Learning*, pages 7279–7289. PMLR, 2021.
- 391 [25] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu.
392 Towards deep learning models resistant to adversarial attacks. In *International Conference on
393 Learning Representations*, 2018.
- 394 [26] Muhammad Muzammal Naseer, Salman H Khan, Muhammad Haris Khan, Fahad Shahbaz Khan,
395 and Fatih Porikli. Cross-domain transferability of adversarial perturbations. *Advances in Neural
396 Information Processing Systems*, 32, 2019.
- 397 [27] Neehar Peri, Neal Gupta, W Ronny Huang, Liam Fowl, Chen Zhu, Soheil Feizi, Tom Goldstein,
398 and John P Dickerson. Deep k-nn defense against clean-label data poisoning attacks. In
399 *European Conference on Computer Vision*, pages 55–70. Springer, 2020.
- 400 [28] Zeyu Qin, Yanbo Fan, Hongyuan Zha, and Baoyuan Wu. Random noise defense against
401 query-based black-box attacks. *Advances in Neural Information Processing Systems*, 34, 2021.
- 402 [29] Evgenia Rusak, Lukas Schott, Roland S Zimmermann, Julian Bitterwolf, Oliver Bringmann,
403 Matthias Bethge, and Wieland Brendel. A simple way to make neural networks robust against
404 diverse image corruptions. In *European Conference on Computer Vision*, pages 53–69. Springer,
405 2020.
- 406 [30] Aniruddha Saha, Akshayvarun Subramanya, and Hamed Pirsiavash. Hidden trigger backdoor
407 attacks, 2019.
- 408 [31] Avi Schwarzschild, Micah Goldblum, Arjun Gupta, John P Dickerson, and Tom Goldstein. Just
409 how toxic is data poisoning? a unified benchmark for backdoor and data poisoning attacks.
410 *arXiv preprint arXiv:2006.12557*, 2020.
- 411 [32] Ali Shafahi, W. Ronny Huang, Mahyar Najibi, Octavian Suci, Christoph Studer, Tudor
412 Dumitras, and Tom Goldstein. Poison frogs! targeted clean-label poisoning attacks on neural
413 networks, 2018.
- 414 [33] Reza Shokri et al. Bypassing backdoor detection algorithms in deep learning. In *2020 IEEE
415 European Symposium on Security and Privacy (EuroS&P)*, pages 175–183. IEEE, 2020.
- 416 [34] Hossein Souri, Micah Goldblum, Liam Fowl, Rama Chellappa, and Tom Goldstein. Sleeper
417 agent: Scalable hidden trigger backdoors for neural networks trained from scratch. *arXiv
418 preprint arXiv:2106.08970*, 2021.

- 419 [35] Jacob Steinhardt, Pang Wei Koh, and Percy Liang. Certified defenses for data poisoning attacks,
420 2017.
- 421 [36] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfel-
422 low, and Rob Fergus. Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199*,
423 2013.
- 424 [37] Lue Tao, Lei Feng, Jinfeng Yi, Sheng-Jun Huang, and Songcan Chen. Better safe than sorry:
425 Preventing delusive adversaries with adversarial training. *Advances in Neural Information*
426 *Processing Systems*, 34, 2021.
- 427 [38] Brandon Tran, Jerry Li, and Aleksander Madry. Spectral signatures in backdoor attacks. In
428 *Advances in Neural Information Processing Systems*, pages 8000–8010, 2018.
- 429 [39] Alexander Turner, Dimitris Tsipras, and Aleksander Madry. Clean-label backdoor attacks.
430 2018.
- 431 [40] Bolun Wang, Yuanshun Yao, Shawn Shan, Huiying Li, Bimal Viswanath, Haitao Zheng, and
432 Ben Y Zhao. Neural cleanse: Identifying and mitigating backdoor attacks in neural networks.
433 In *2019 IEEE Symposium on Security and Privacy (SP)*, pages 707–723. IEEE, 2019.
- 434 [41] Maurice Weber, Xiaojun Xu, Bojan Karlaš, Ce Zhang, and Bo Li. Rab: Provable robustness
435 against backdoor attacks. *arXiv preprint arXiv:2003.08904*, 2020.
- 436 [42] Cihang Xie, Zhishuai Zhang, Yuyin Zhou, Song Bai, Jianyu Wang, Zhou Ren, and Alan L Yuille.
437 Improving transferability of adversarial examples with input diversity. In *Proceedings of the*
438 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2730–2739, 2019.
- 439 [43] Chen Zhu, W Ronny Huang, Hengduo Li, Gavin Taylor, Christoph Studer, and Tom Goldstein.
440 Transferable clean-label poisoning attacks on deep neural nets. In *International Conference on*
441 *Machine Learning*, pages 7614–7623, 2019.

442 Checklist

- 443 1. For all authors...
- 444 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
445 contributions and scope? [Yes]
- 446 (b) Did you describe the limitations of your work? [Yes]
- 447 (c) Did you discuss any potential negative societal impacts of your work? [N/A]
- 448 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
449 them? [Yes]
- 450 2. If you are including theoretical results...
- 451 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 452 (b) Did you include complete proofs of all theoretical results? [N/A]
- 453 3. If you ran experiments...
- 454 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
455 mental results (either in the supplemental material or as a URL)? [Yes]
- 456 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
457 were chosen)? [Yes]
- 458 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
459 ments multiple times)? [Yes]
- 460 (d) Did you include the total amount of compute and the type of resources used (e.g., type
461 of GPUs, internal cluster, or cloud provider)? [Yes]
- 462 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 463 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 464 (b) Did you mention the license of the assets? [N/A]
- 465 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]

- 466 (d) Did you discuss whether and how consent was obtained from people whose data you're
467 using/curating? [N/A]
- 468 (e) Did you discuss whether the data you are using/curating contains personally identifiable
469 information or offensive content? [N/A]
- 470 5. If you used crowdsourcing or conducted research with human subjects...
- 471 (a) Did you include the full text of instructions given to participants and screenshots, if
472 applicable? [N/A]
- 473 (b) Did you describe any potential participant risks, with links to Institutional Review
474 Board (IRB) approvals, if applicable? [N/A]
- 475 (c) Did you include the estimated hourly wage paid to participants and the total amount
476 spent on participant compensation? [N/A]