

GYM: A Comprehensive Defense Approach against DNN Backdoor Attacks

Yi Zeng¹, Han Qiu², Shangwei Guo³, Tianwei Zhang³, Meikang Qiu⁴, Bhavani Thuraisingham⁵

²Telecom Paris, Institut Polytechnique de Paris, Palaiseau, France, 91120.

¹University of California San Diego, CA, USA, 92122.

³Nanyang Technological University, Singapore, 639798.

⁴Columbia University, New York, USA, 10027.

⁵The University of Texas at Dallas, Texas, USA, 75080.

Abstract

Public resources and services (e.g., datasets, training platforms, pre-trained models) have been widely adopted to ease the development of deep learning-based applications. However, if the third-party providers are untrusted, they can inject poisoned samples into the datasets or embed backdoors in those models. Such an integrity breach can cause severe consequences, especially in safety- and security-critical applications. Various backdoor attack techniques have been proposed for higher effectiveness and stealthiness. Unfortunately, existing defense solutions are not practical to thwart those attacks in a comprehensive way.

In this paper, we propose GYM, a novel and effective defense solution to defeat different types of backdoor attacks and enhance DL models' robustness. The key innovations of our approach are two preprocessing functions: (1) an intensive function is used to transform clean images for fine-tuning of the infected model. This can invalidate the effects of embedded backdoors; (2) a lightweight function is adopted to invalidate triggers during inference. The combination of these two functions in two stages can achieve reliable and comprehensive protection of backdoored models. Extensive experiments show that our solution can effectively mitigate six different kinds of backdoor attacks and outperform four state-of-the-art defense solutions for various DNN models and datasets.

Introduction

The past several years have witnessed the rapid development of Deep Learning (DL) technology. Various DL models today are widely adopted in many scenarios, e.g., image classification (Chan et al. 2015), speech recognition (Deng and Platt 2014), natural language processing (Collobert and Weston 2008). These applications significantly enhance the quality of life and work efficiency. With the increased complexity of Artificial Intelligence tasks, more sophisticated DL models need to be trained, which require large-scale datasets and a huge amount of computing resources.

To reduce the training cost and effort, it is now common for developers to leverage third-party resources and services for efficient model training. Developers can download state-of-the-art models from the public model zoos or purchase them from model vendors. They can also download or purchase valuable datasets from third parties and train the models by themselves. A more convenient way is to utilize public cloud services (e.g., Amazon SageMaker (Liberty et al.

2020), GoogleVision AI (Hosseini, Xiao, and Poovendran 2017), Microsoft Computer Vision (Han et al. 2013), etc.), which can automatically deploy the training environment and allocate hardware resources based on users' demands.

However, new security threats are introduced to DNN models when the third party is not trusted. One of the most severe threats is the DNN backdoor attacks (Li et al. 2020a): the adversary injects a backdoor into the victim model, causing it to behave normally over benign samples, but predict the samples with an attacker-specified trigger as wrong labels desired by the adversary. Typically, a backdoor injection can be achieved by directly modifying the neurons (Liu et al. 2017) or poisoning the training datasets (Gu, Dolan-Gavitt, and Garg 2017). In practice, the developer may obtain a poisoned dataset if the source is untrusted. It is hard to detect such a threat as a very small ratio of malicious samples can lead to a backdoored model. When the developer outsources the model training task to an untrusted cloud provider, the adversary can inject the backdoor by either dataset poisoning or parameter modifications. It will then be difficult for the developer to detect the existence of backdoors, as the model only has anomalous predictions on samples with triggers, which are agnostic to the developer.

It is of paramount importance to have an effective method to address these severe threats. Past works proposed some approaches to detect the existence of backdoors or eliminate them from the infected models. Unfortunately, most of them have certain limitations. First, some solutions require the defender to have poisoned data samples (Chen et al. 2018; Du, Jia, and Song 2019), or knowledge of the attack techniques (Xu et al. 2019) and triggers (Chou et al. 2018). This assumption is not held when the defender is only given the victim model. Second, these solutions are not comprehensive to cover all different types of attack techniques and trigger patterns. For instance, (Wang et al. 2019) is effective against the single target attack but fails to identify the all-to-all attack where there is more than one target label for the malicious samples (Gu, Dolan-Gavitt, and Garg 2017). We will empirically validate this in our evaluation section. (Liu et al. 2019) cannot defeat attacks with complex trigger patterns (e.g., watermarking in the background), as claimed in that paper. More importantly, most of these defense works only consider traditional backdoor attacks, while ignoring the recently discovered advanced attacks (e.g., invisible backdoor

88 attacks (Li et al. 2019)). As one of our contributions, we will
89 categorize past defense solutions and analyze their limita-
90 tions in this paper.

91 Motivated by the gap between the severity of backdoor at-
92 tacks and the limitations of existing solutions, we propose
93 GYM, a novel and comprehensive method to defeat various
94 DNN backdoor attacks. A successful backdoor attack relies
95 on both the backdoor in the infected model and effective trig-
96 gers hidden in the malicious samples. Thus, the key idea of
97 our approach is the integration of model fine-tuning (which
98 is used to weaken the effects of the backdoor) and input pre-
99 processing (which is used to affect the impact of triggers).
100 Given an infected model, our solution has two steps. During
101 the fine-tuning phase, it retrains the model for a few epochs
102 using some intensive preprocessed clean data samples¹. Dur-
103 ing the inference phase, each data sample (either clean one or
104 trigger-patched one) is first transformed by a lightweight pre-
105 processing function and then fed into the fine-tuned model
106 for prediction. With these two steps, the model will correct
107 the labels of malicious samples while maintaining high per-
108 formance for normal data.

109 Our method is effective against various backdoor attacks
110 and trigger patterns and does not need any prior knowledge
111 about the attack techniques or poisoned samples. We conduct
112 comprehensive evaluations to validate our solution. We con-
113 sider different datasets (Cifar10, GTSRB, PubFig) and mod-
114 els (ResNet-18, LeNet-8, VGG-16). We implement different
115 attack techniques (BadNet (Gu, Dolan-Gavitt, and Garg
116 2017), Neural Trojan (Liu et al. 2017), invisible backdoor
117 (Li et al. 2019)), different triggers (Square, watermark, ad-
118 versarial perturbation) and different modes (single target, all-
119 to-all). Our method can successfully defeat all these threats.
120 Evaluations also show that our method can outperform four
121 state-of-the-art works (Neural Cleanse (Wang et al. 2019),
122 Fine-pruning (Liu, Dolan-Gavitt, and Garg 2018), FLIP, and
123 ShrinkPad-4 (Li et al. 2020b)) in defeating different back-
124 door attacks and enhancing the model robustness.

125 Background about Backdoor Attacks

126 Given a DNN model f_θ with parameters θ , a backdoor at-
127 tack can be formulated as a tuple $(\Delta\theta, \delta)$, where $\Delta\theta$ is the
128 backdoor injected by the adversary to the model parameters,
129 and δ is an attacker-specified trigger. Then the compromised
130 model $f_{\theta+\Delta\theta}$ still has state-of-the-art performance for nor-
131 mal samples: $f_{\theta+\Delta\theta}(x) = f_\theta(x), \forall x \in \mathcal{X}$. However, for an
132 input sample containing the trigger, the model will predict a
133 label different from the correct one: $y' = f_{\theta+\Delta\theta}(x + \delta) \neq$
134 $f_{\theta+\Delta\theta}(x), \forall x \in \mathcal{X}$. y' can be a fixed label pre-determined
135 by the attacker, or an arbitrary unmatched label.

136 The adversary has multiple ways to embed the backdoor
137 into the DNN model. (1) Data poisoning (Gu, Dolan-Gavitt,
138 and Garg 2017; Chen et al. 2017): the adversary generates a
139 number of poisoned samples with the desired labels and in-
140 corporates such samples into the clean training set to train a
141 backdoor model. (2) Direct modification (Liu et al. 2017): the

¹This step can be applied to the case where the training set is
poisoned as well. The defender can use this intensive function to
preprocess the dataset and then trains the model from scratch.

adversary can select critical neurons and weights for modifi- 142
cation via model retraining. (3) Transfer learning (Yao et al. 143
2019): if the adversary injects backdoor into a teacher model, 144
the student models transferred from this teacher model may 145
still contain the backdoor. 146

There can be different designs for malicious triggers. The 147
most common one is a small block with several pixels placed 148
at the corner of the image. For instance, (Gu, Dolan-Gavitt, 149
and Garg 2017) added a white square onto the right bottom 150
of the image as the trigger. (Liu et al. 2017) introduced a 151
colored square to activate the backdoor. 152

It is possible that the trigger size is big and located across 153
the images. Such patterns need to be designed not to affect 154
the clean samples. For instance, watermarks are embedded 155
over the background of the samples (Liu et al. 2017; Chen 156
et al. 2017). A special pair of glasses function as a trigger 157
when it is worn by a person (Chen et al. 2017). 158

The third type is invisible triggers, introduced in (Liao 159
et al. 2018; Li et al. 2019). Inspired by the adversarial exam- 160
ples, such triggers are imperceptible perturbations, which are 161
visually indistinguishable from normal samples. These trig- 162
gers can make the corresponding backdoor attacks stealthier, 163
and it is hard to detect poisoned data from the training set. 164

165 Defense Requirements and Existing Solutions

166 Defense Requirements

To effectively defeat DNN backdoor attacks, a good solution 167
must have the following properties. 168

- 169 • *Robust*: the solution should be capable of effectively de- 170
tecting or eliminating backdoor with a low attack success 171
rate. It should be hard to evade this solution even if the 172
adversary knows the defense mechanism. 173
- 174 • *Attack-agnostic*: the defender does not have any knowl- 175
edge of the employed attack technique, trigger informa- 176
tion (pattern, location, size, desired label, etc.). He does 177
not have access to the poisoned data samples either. All he 178
has are clean data samples and a suspicious model that can 179
be potentially infected with backdoors. 180
- 181 • *Comprehensive*: the defense solution should be able to 182
cover different types of backdoor attacks, regardless of the 183
size, complexity, and visibility of triggers, as well as the 184
attacker’s target labels. 185
- 186 • *Functionality-preserving*: this solution should have a 187
small impact on the model performance of clean samples. 188
- 189 • *Lightweight*: the defender should be able to defeat back- 190
door attacks in a lightweight manner. Given a suspicious 191
model, the defense cost should be much smaller than train- 192
ing a clean model from scratch. During inference, the pre- 193
diction process cannot incur high overhead, either. 194

190 Review of Existing Solutions

Various defense techniques against backdoor attacks have 191
been proposed. We classify them into different categories and 192
check their satisfaction of the above requirements. 193

Backdoor Detection. The most popular direction is to check 194
if one DL model has a backdoor injected. (Wang et al. 2019) 195

196 made an attempt towards this goal with boundary outlier de- 255
197 tection. Some works followed a similar idea to detect the ex- 256
198 istence of backdoors and adopted different techniques to re- 257
199 cover the trigger, such as Generative Adversarial Networks 258
200 (Chen et al. 2019), new regularization terms (Guo et al.
201 2019), Generative Distribution Modeling (Qiao, Yang, and
202 Li 2019), and Artificial Brain Stimulation (Liu et al. 2019).

203 Unfortunately, these approaches make two unrealistic as- 260
204 sumptions. First, they assume there is only one target label 261
205 for all malicious samples (i.e., single-target attack). They 262
206 are ineffective when the adversary has more than one tar- 263
207 get labels (e.g., all-to-all attack (Gu, Dolan-Gavitt, and Garg 264
208 2017)). Second, they assume the trigger must have a small 265
209 size and simple pattern. They fail to detect complex triggers 266
210 such as watermarks in the background. Hence, these solu- 267
211 tions cannot meet the *comprehensiveness* requirement. 268

212 (Xu et al. 2019) proposed another detection approach 269
213 without the above assumptions. It tries to build a classifier 270
214 to distinguish benign and infected models. It needs to mimic 271
215 all possible backdoor attacks to build all those models, which 272
216 is costly and impractical as there are too many existing and 273
217 unknown ways to perform backdoor attacks on DL models. 274
218 This solution is thus not *lightweight*. 275

219 **Backdoor Invalidation.** This direction is to remove the po- 276
220 tential backdoor from the model directly without detection. 277
221 (Liu, Dolan-Gavitt, and Garg 2018) proposed to use fine- 278
222 pruning and fine-tuning to break the backdoor effects. How- 279
223 ever, this solution may reduce the prediction accuracy over 280
224 clean samples, which is not *functionality-preserving*. 281

225 **Trigger Detection.** Instead of checking the suspicious 282
226 model, this direction focuses on the samples with triggers. 283
227 It can be applied to two cases. The first case is to detect 284
228 if the training data set contains poisoned samples. For in- 285
229 stance, (Chen et al. 2018) discovered that normal and poi- 286
230 soned data yield different features in the last hidden layer’s 287
231 activations. (Tran, Li, and Madry 2018) proposed a new rep- 288
232 resentation to classify benign and malicious samples. (Du, 289
233 Jia, and Song 2019) adopted differential privacy to detect ab- 290
234 normal training samples. These solutions cannot work when 291
235 the defender only has the infected model rather than the poi- 292
236 soned data samples, especially when the backdoor is injected 293
237 via direct neuron modification instead of data poisoning (Liu 294
238 et al. 2017). They cannot achieve *comprehensiveness*. 295

239 The second case is the online detection of triggers in the 296
240 inference samples. (Gao et al. 2019) proposed to superim- 297
241 pose a target sample with a benign one from a different class. 298
242 A benign sample’s prediction result will be altered while a 299
243 malicious sample will still keep the same due to the triggers. 300
244 However, this approach may not be *robust* when the super- 301
245 imposed benign image has overlap with the trigger. (Chou 302
246 et al. 2018) proposed to use image processing techniques 303
247 (e.g., Grad-CAM) to visualize and reveal the trigger. This ap- 304
248 proach requires prior knowledge of the trigger pattern, which 305
249 is not *attack-agnostic*. 306

250 **Trigger Invalidation.** The last direction is to directly inval- 307
251 idate the effects of the triggers from the inference samples. 308
252 (Li et al. 2020b) proposed to adopt image preprocessing to 309
253 transform input such that the backdoor model will give cor- 310
254 rect results for both benign and malicious samples. However,

since backdoor models and triggers have much higher robust- 255
ness than adversarial attacks, this solution is not *comprehen- 256*
sive, as it can only handle simple triggers, but fail to defeat 257
complex ones (e.g., watermarks in the background). 258

Our Proposed Method 259

In this section, we present the details of our proposed solu- 260
tion, GYM, to defeat DNN backdoor attacks. As discussed in 261
the previous section, it is difficult to counter backdoor at- 262
tacks and meet those requirements using just one defense 263
direction. So our method will *combine the mechanisms of 264*
both backdoor invalidation and trigger invalidation. Figure 265
1 illustrates the methodology overview. It consists of two 266
stages. The first stage is fine-tuning: we introduce an inten- 267
sive preprocessing function to transform a small number of 268
clean data samples, which will be used to fine-tune the in- 269
fected model. The second stage is inference: we design a 270
lightweight function to preprocess the inference samples, and 271
then send the transformed output to the fine-tuned model for 272
prediction. This function can remove the effects of triggers 273
while still preserving the model’s performance over clean 274
samples. Below we describe each step in detail. 275

Stage 1: Fine-tuning with Intensive Preprocessing 276

We fine-tune the model with preprocessed data samples to 277
weaken the malicious impact of injected backdoors. Our 278
method is different from previous works that directly fine- 279
tune or fine-prune the model (Liu, Dolan-Gavitt, and Garg 280
2018; Wang et al. 2019), as they can compromise the model 281
usability. We prepare a fine-tuning dataset with a preprocess- 282
ing function consisting of some transformations (Figure 2): 283

T1: Optical distortion. The optical distortion used here is 284
a pincushion distortion (Liu, Malcolm, and Xu 2010), where 285
image magnification increases with the distance from the op- 286
tical axis. In Figure 2, we can observe that lines that do not go 287
through the center of the image are bowed towards the center 288
after this transformation, like a pincushion. We conduct this 289
procedure to map the representation of inputs away from the 290
original representation in the hyperdimensional space. As the 291
accuracy of clean data can be recovered with our fine-tuning, 292
malicious samples during inference will have a lower success 293
rate with the shifted decision boundary. 294

T2: Three median filters in different spaces. A set of three 295
median filters are employed to reinforce the model against 296
the backdoors by fine-tuning with the preprocessed samples. 297

The first median filter is performed in the gamma space 298
with a gamma compression. This filter can help the model ac- 299
quire adaptation against strong perturbation cost by a median 300
filter. We set the encoding gamma value as 0.6 to lighten the 301
images. This gamma compression causes large-value pixels 302
inside the image to bend in together. The small-value pixels 303
in the image thus have a better contrast against large-value 304
pixels. Therefore, the median filter can better smoothen those 305
pixels. The kernel size of the median filter is 5×5 . 306

The second median filter is also performed in the gamma 307
space but with a gamma extension. We first multiply each 308
pixel with the multiplier (set as 1.53) to further lighten up 309
the images to bend large-value pixels together and disrupt the 310

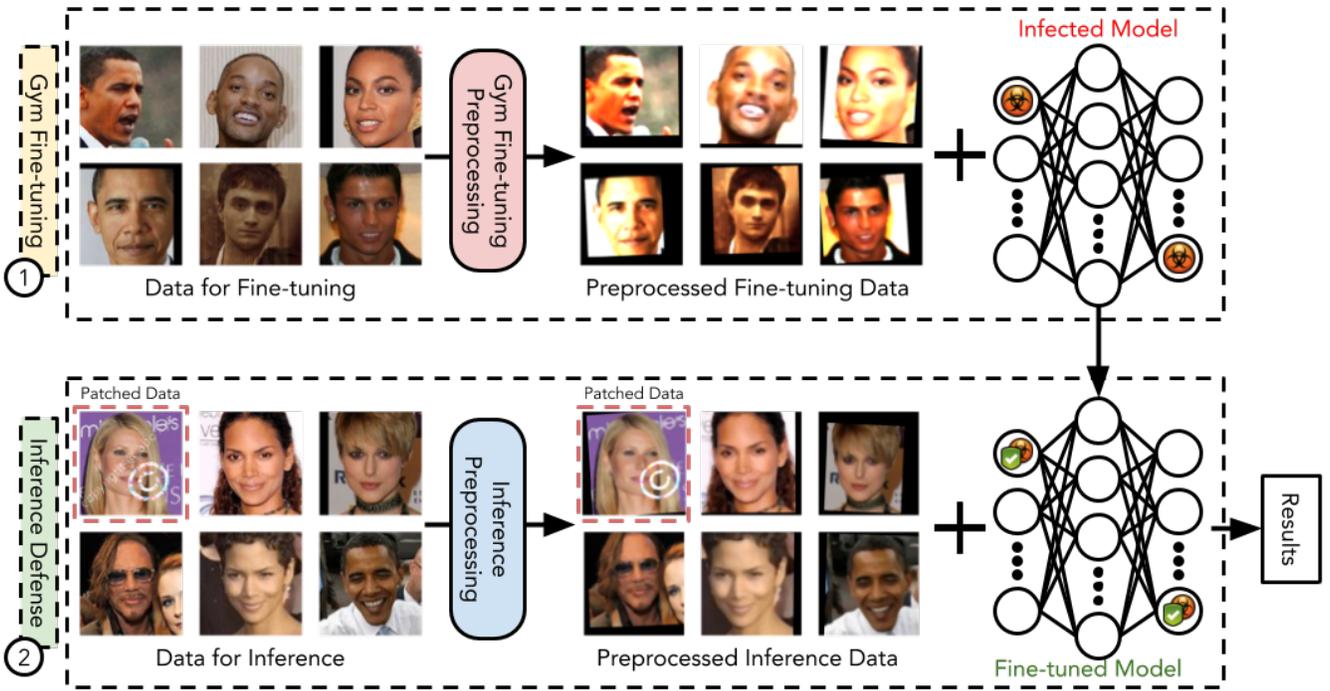


Figure 1: The workflow of GYM to defeat DNN backdoor attacks (using PubFig dataset as an example): ① An intensive preprocessing function is introduced over the inputs, the output of which is used to fine-tune the infected model. ② The inference samples are transformed via an inference preprocessing function to invalidate triggers for the fine-tuned model.

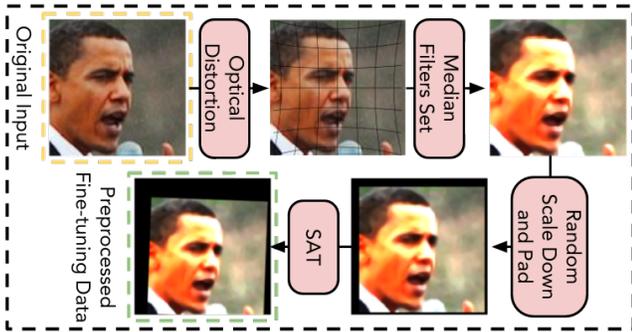


Figure 2: The intensive preprocessing function includes three affine transformations (optical distortion, random scale down with padding, and SAT (Zeng et al. 2020)) and one set of three median filters in different spaces.

As shown in Figure 2, with three median filters, a per- 322
 turbed result is obtained. The entire transformation can 323
 strengthen the infected model and shift the decision bound- 324
 ary away with preprocessed data. This guarantees the model 325
 is robust over clean samples. 326

T3: Random scaling down with padding. This procedure 327
 scales down the image and then resizes it up in a 328
 random manner. Inspired by ShinkPad (Li et al. 2020b), which 329
 demonstrated the capability of invalidating BadNets (Gu, 330
 Dolan-Gavitt, and Garg 2017) triggers, our transformation 331
 improves the randomness level. First, we randomly scale the 332
 image into a smaller size ranging between $[0.8-1]$ of the orig- 333
 inal size, by dropping random pixels. We then pad it to the 334
 original size by randomly choosing a point as the center. 335
 Such operation can shift all the pixels away from the ac- 336
 tual coordinates. Thus, samples will likely move away from 337
 the infected model’s original representation output (with a 338
 certain accuracy drop). This procedure can harden the DL 339
 models’ boundary to help the fine-tuned model adapt to the 340
 pixels shifting around. Thus, when adopting a shifting strat- 341
 egy for preprocessing input samples during the inference, the 342
 model’s accuracy on clean samples can be maintained. 343

T4: Stochastic Affine Transformation (SAT). Finally, we 344
 adopt a preprocessing function in (Zeng et al. 2020) to distort 345
 the image with rotation, scaling, and shifting. SAT first ran- 346
 domly shifts all the pixels horizontally and vertically. Then, 347
 it randomly rotates the image to a certain scale. Finally, it 348
 randomly scales the image up or down to produce the final 349
 output. The visual effect is shown in Figure 2. Note that 350

311 continuity between pixels. Then, a gamma extension is used
 312 to dim the image for obtaining a higher contrast. Third, we
 313 scale down the image to 75% of its original size and conduct
 314 the same median filter as the first one in this gamma exten-
 315 sion space. Such an operation can help the fixed-kernel median
 316 filter remove more outliers globally and obtain smoother
 317 results. Finally, we resize the image back to its original size.

318 The third median filter works with another scaling down
 319 (resize) procedure. We scale down the image to 0.8 of its
 320 original size. The third median filter further smoothens the
 321 pixels in this downscale space.

351 some steps in SAT are similar to the previous transforma-
 352 tion. Adopting these two random transformations can make
 353 the fine-tuned model better adapt to the affine transforma-
 354 tions used during the inference phase.

355 **Fine-tuning the suspicious model.** As shown in Figure 2,
 356 the intensive preprocessing function can introduce significant
 357 distortion to the samples. The performance of the model will
 358 drop over the preprocessed samples. Our solution fine-tunes
 359 the model with the preprocessed samples to help the model
 360 recognize such transformations. GYM only requires a small
 361 number of epochs with a few fine-tuning samples to reach
 362 state-of-the-art performance. Then, the classification bound-
 363 ary of the infected model will be altered against malicious
 364 samples patched with the triggers.

365 Stage 2: Inference with Lightweight Preprocessing

366 In this stage, in order to reduce the computation complex-
 367 ity and overhead of inference, we design a lightweight ver-
 368 sion of the preprocessing function to transform the input be-
 369 fore sending it to the fine-tuned model. This inference pre-
 370 cessing function only includes a set of two median filters
 371 and one affine transformation, as shown in Figure 3. The first
 372 median filter is used to smoothen the pixels in the raw input.
 373 In contrast, the second median filter is integrated with
 374 the scaling down mechanism (i.e., same as the third median
 375 filter in the intensive preprocessing). Finally, the Stochastic
 376 Affine Transformation (SAT) is adopted over the filtered data
 377 to map the pixels away from the original coordinates. With
 378 such transformations, the model can still recognize the clean
 379 samples correctly, while the fine-tuned backdoor cannot rec-
 380 ognize the preprocessed triggers anymore. This makes our
 381 solution robust against both normal and malicious samples.

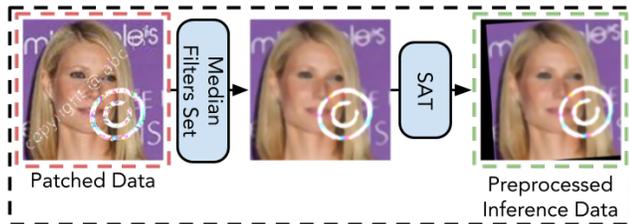


Figure 3: Inference Preprocessing consists of two transforma-
 tions: a set of two median filters affects the triggers from
 two spaces; SAT helps distort the image.

382 Security Analysis

383 We visually interpret the mechanism of our defense solution,
 384 as shown in Figure 4. The left figure shows the infected de-
 385 cision boundary between the original class (red region) and
 386 the attacker-desired class (white region) in a backdoor attack.
 387 We can see such an infected model can still perform well
 388 on clean samples. However, the trigger-patched data in the
 389 source class will be classified to be the target class as the trig-
 390 ger moves them across the boundary. Since the patched sam-
 391 ples still contain features similar to the samples in the source
 392 class, and the trigger impact is small, those patched samples
 393 will be close to the boundary in the hyper-dimensional space.

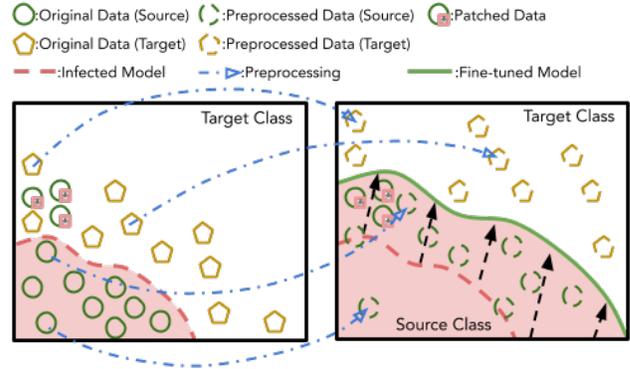


Figure 4: The visual interpretation of GYM over an infected
 model’s decision boundary.

The blue arrows represent our preprocessing function, which
 maps all the clean data used for fine-tuning away from the
 original location. The fine-tuning process can shift the in-
 fected decision boundary adaptive to the fine-tuning data.
 Since patched data are not used in fine-tuning, their hyperdi-
 mensional representation will not be shifted with the others.
 Thus, the new boundary can correct the prediction of patched
 data into the source class again.

Evaluation

Implementation

We conduct a comprehensive evaluation of our proposed de-
 fense against different kinds of backdoor attacks. Table 1
 summarizes the configurations of these attacks, as well as
 the target models and datasets. We mainly replicate the same
 implementations as the original attack papers.

Dataset	Model	Attack	Target Label	Poisoning Ratio
Cifar10	ResNet-18	Trojan (WM)	'7'	10%
		Trojan (SQ)	'7'	10%
		BadNets All-to-all	'i+1'	10%
		L2 invisible	'3'	5%
		L0 invisible	'4'	5%
GTSRB	LeNet-8	BadNets	'33'	10%
PubFig	VGG-16	Trojan (WM)	'0'	10%
		Trojan (SQ)	'0'	10%

Table 1: Datasets and backdoor attacks in evaluation.

Cifar10 is a widely-adopted dataset for image classifica-
 tion. It contains 50000 training images and 10000 testing im-
 ages. We adopt ResNet-18 (He et al. 2016) to train five back-
 door models with different triggers. We inject 10% of poi-
 soned samples into the training set to generate the first three
 models, while the last two models have a poisoning ratio of
 5%, which is enough for the invisible backdoor attacks. All
 the infected models are trained with Adadelata (Zeiler 2012)
 as the optimizer and an initial learning rate of 0.05 for 200
 epochs. Specifically, for the first two models, we implement
 the trojan attack in (Liu et al. 2017) with the trigger of the wa-
 termark (WM) and square (SQ), respectively. The attacker’s

Model	Attack	Baseline			Inference (I)		Fine-tuning + Inference (I)		Fine-tuning + Inference (L)	
		ACC	ASR	TCP	ACC	ASR	ACC	ASR	ACC	ASR
ResNet-18 (Cifar10)	Trojan (WM)	0.830	1.000	0.075	0.520	0.810	0.805	0.130	0.785	0.045
	Trojan (SQ)	0.880	1.000	0.100	0.600	0.635	0.760	0.065	0.780	0.040
	BadNets All-to-all	0.875	0.670	0.125	0.435	0.150	0.765	0.020	0.670	0.030
	L2 invisible	0.900	0.985	0.110	0.610	0.420	0.790	0.205	0.810	0.180
	L0 invisible	0.895	0.990	0.070	0.645	0.135	0.805	0.080	0.825	0.080
LeNet-8 (GTSRB)	BadNets	0.960	0.985	0.020	0.660	0.170	0.875	0.045	0.905	0.035
VGG-16 (PubFig)	Trojan (WM)	0.960	1.000	0.025	0.400	0.360	0.840	0.010	0.910	0.010
	Trojan (SQ)	0.955	1.000	0.015	0.400	0.055	0.815	0.015	0.870	0.015

Table 2: Evaluation of ACC and ASR with different techniques of GYM: (I) denotes intensive preprocessing, while (L) denotes lightweight preprocessing.

	ResNet-18 (Cifar10)									
	Trojan (WM)		Trojan (SQ)		BadNets (All-to-all)		L2 invisible		L0 invisible	
	ACC	ASR	ACC	ASR	ACC	ASR	ACC	ASR	ACC	ASR
No Defense	0.830	1.000	0.880	1.000	0.875	0.670	0.900	0.985	0.895	0.990
GYM	0.785	0.045	0.780	0.040	0.765	0.020	0.810	0.180	0.825	0.080
Neural Cleanse (unlearning)	0.895	0.085	0.910	0.155	NA	NA	NA	NA	NA	NA
Fine-pruning	0.835	0.195	0.845	0.235	0.630	0.055	0.860	0.990	0.860	0.880
Fine-pruning (finetuned)	0.855	0.650	0.870	0.140	0.775	0.055	0.895	0.935	0.900	0.810
FLIP	0.830	0.880	0.775	0.090	0.855	0.020	0.900	0.965	0.890	0.975
ShrinkPad-4	0.720	1.000	0.800	0.075	0.625	0.130	0.855	0.735	0.850	0.985

Table 3: Comparison of ACC and ASR between GYM with previous defense methods for models on the Cifar10 dataset.

target label is set as class ‘7:Horse’. For the third model, we replicate the all-to-all attack in BadNets (Gu, Dolan-Gavitt, and Garg 2017): the trigger is a white square of 5×5 pixels located at the right bottom of the image. The target label of a sample from class i is set to be class $i + 1$. For the last two ResNet-18 models, we replicate the L2 and L0 invisible attacks in (Li et al. 2019). The target class is obtained by forward-passing the trigger to a pre-trained clean ResNet-18 model: ‘3:Cat’ for L2 attack and ‘4:Deer’ for L0 attack.

The second dataset included in the evaluation is the GTSRB (Stallkamp et al. 2012), which contains 35228 training samples and 12630 testing samples in 43 classes. Different from the Cifar10 case, we directly obtain a backdoor model (LeNet-8) from (Wang et al. 2019), which has been compromised by the BadNets technique (Gu, Dolan-Gavitt, and Garg 2017). The trigger is a white square of 5×5 pixels at the right bottom, and the target label is ‘33:turn right ahead’.

The last dataset for evaluation is the PubFig (Kumar et al. 2009), which contains 11070 training images and 2768 testing images of 83 celebrities. We also direct get two backdoor models (VGG-16) from (Jin et al. 2020), compromised by the Trojan attacks (Liu et al. 2017). The trigger patterns are the same as the Cifar10 case (WM and SQ) and the target label is ‘0:Adam Sandler’.

We use Keras with Tensorflow backend as the DL framework for the implementations. We adopt model accuracy (ACC) over clean samples and attack success rate (ASR) over patched samples to quantify a DL model’s robustness. ACC is measured using 200 clean samples, and ASR is calculated using 200 different samples patched with the corresponding triggers. A good defense should be effective (low ASR) and functionality-preserving (high ACC). We conduct

all the experiments on a server equipped with 8 Intel I7-7700k CPUs and 4 NVIDIA GeForce GTX 1080 Ti GPUs.

Effectiveness of Each Technique

We first evaluate the ACC and ASR with different techniques from GYM. The results are summarized in Table 2. The Baseline column shows the results of the infected model without any defense. We use a new metric, Target Class Probability (TCP), to denote the percentage that the infected model will predict the clean testing data as the target class. Clearly, a good defense should make the ASR close to or even below TCP, in order to defeat the attacks. The Inference (I) column demonstrates the results when we just use intensive preprocessing for inference without fine-tuning the model. The last two columns in Table 2 present the results when we fine-tune the model with intensive preprocessing and then perform inference with intensive preprocessing and lightweight preprocessing, respectively.

For the inference with intensive preprocessing, we can observe a severe drop in ACC due to the absence of fine-tuning. However, the ASR of only four infected models drops below 0.2. This indicates that solely adopting a strong preprocessing function cannot effectively invalidate the trigger, as many state-of-the-art backdoor triggers are preprocessing-robust (e.g., Trojan WM, L2 invisible). A common trait of those kinds of robust triggers is that the trigger size is close to the sample space, making them hard to be removed as almost all the pixels in the patched data will be affected.

For the Fine-tuning+Inference (I), we use the intensive function to preprocess 10000 clean samples and then use them to fine-tune the model for five epochs. Compared to the cost of training a new model from scratch, which can

421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452

453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483

484 take more than 200 epochs using 50000 training samples in
 485 our case, our fine-tuning is very lightweight. As shown in
 486 Table 2, we observe that the fine-tuned model can increase
 487 the ACC and significantly decrease ASR on intensive pre-
 488 processed data, even for preprocessing-robust triggers. This
 489 conforms to the security analysis in the previous section that
 490 the decision boundary can shift with the fine-tuning data and
 491 encompass the patched data as the ACC rises. It is worth not-
 492 ing that the ASR of the L2 invisible attack is still at a rela-
 493 tively high level (0.205). This calls for a stronger preprocess-
 494 ing mechanism to tackle such a stealthy attack.

495 Finally, we check our complete solution, where intensive
 496 preprocessing is adopted for fine-tuning, and lightweight pre-
 497 processing is used for inference. As shown in the column of
 498 Fine-tuning + Inference (L) in Table 2, we can conclude that,
 499 in most cases, the lightweight preprocessing can further in-
 500 crease the ACC and reduce ASR. One counterexample here
 501 is the BadNets All-to-all attack. where ACC has a relative
 502 large drop. The reason is that the patch data of this infected
 503 model cover all the decision boundaries across all the classes.
 504 Thus, reducing the defense scale can cause a worse result
 505 than inference with the intensive preprocessing. This case
 506 can be resolved by adding more data to the fine-tuning set
 507 and increasing the number of tuning epochs. Our solution
 508 works better for the attacks with fewer target labels. In a nut-
 509 shell, the proposed defense can successfully mitigate all the
 510 attacks in our consideration.

511 Comparison with Existing Defense Methods

512 Next, we compare GYM with some existing solutions: Neural
 513 Cleanse with Unlearning (Wang et al. 2019), Fine-pruning
 514 (Liu, Dolan-Gavitt, and Garg 2018), FLIP, and ShrinkPad-4
 515 (Li et al. 2020b). To have a fair comparison, for all defense
 516 methods based on fine-tuning, we set the number of available
 517 clean samples as 10000. For Fine-pruning, we only prune the
 518 last convolutional layer of the infected model. We stop the
 519 pruning process when the validation accuracy is decreased
 520 by 4% compared to the baseline ACC, as suggested in (Liu,
 521 Dolan-Gavitt, and Garg 2018). The number of epochs to fine-
 522 tune the pruned model is one.

523 Table 3 shows the comparison results using the Cifar10
 524 dataset, where we train the backdoor model from differ-
 525 ent poisoned datasets. We can observe that GYM gets the
 526 best defense results than other solutions. Particularly, Neu-
 527 ral Cleanse fails to detect the backdoor caused by the Bad-
 528 Nets All-to-all technique as it assumes there is only one tar-
 529 get label. Also, since the original design of Neural Cleanse
 530 did not consider invisible attacks, its out layer detector can-
 531 not discern the target class. Hence, it fails to detect invis-
 532 ible backdoor attacks as well². FLIP and ShrinkPad-4 are
 533 not able to tackle complex triggers such as watermarks or
 534 imperceptible perturbations. This confirms the limitations of
 535 preprocessing-only approaches.

536 Table 4 and 5 present the comparisons for the backdoor
 537 models on GTSRB and PubFig datasets. The infected mod-

²We can manually set the target label in the detector to make it work. However, this is impractical due to the violation of attack-agnostic requirement

	LeNet-8 (GTSRB)	
	BadNets	
	ACC	ASR
No Defense	0.960	0.985
GYM	0.905	0.035
Neural Cleanse (unlearning)	0.960	0.190
Fine-pruning	0.930	0.020
Fine-pruning (finetuned)	0.940	0.545
FLIP	0.535	0.005
ShrinkPad-4	0.945	0.080

Table 4: Comparison of ACC and ASR with past defenses for the model on the GTSRB dataset.

	VGG-16 (PubFig)			
	Trojan (WM)		Trojan (SO)	
	ACC	ASR	ACC	ASR
No Defense	0.960	1.000	0.955	1.000
GYM	0.910	0.010	0.870	0.015
Neural Cleanse (unlearning)	0.880	0.025	0.810	0.010
Fine-pruning	0.909	1.000	0.855	1.000
Fine-pruning (finetuned)	0.929	1.000	0.895	1.000
FLIP	0.930	0.385	0.915	0.015
ShrinkPad-4	0.960	0.995	0.940	0.015

Table 5: Comparison of ACC and ASR with past defenses for models on the PubFig dataset.

538 els are downloaded directly online. We can also observe the
 539 advantages of GYM over other solutions. It is worth not-
 540 ing that the ASR of Fine-pruning maintains 1 on the PubFig
 541 dataset, indicating that a fixed early stop criterion of 4% ac-
 542 curacy drop in ACC is not effective and generalizable. The
 543 defender cannot monitor the ASR to determine the optimal
 544 moment to stop the fine-pruning and balance the security-
 545 usability trade-off. Hence, it is impractical to apply this tech-
 546 nique when the defender is attack-agnostic. In contrast, GYM
 547 can be used easily without this concern.

548 Conclusion

549 This paper proposes GYM, a novel and efficient solution to
 550 mitigate backdoor attacks against DL models. Unlike past
 551 works focusing on either infected models or triggers, our so-
 552 lution adopts novel techniques to break the effects of both
 553 backdoors in the models and triggers in the input samples.
 554 We first design a novel fine-tuning technique with intensive
 555 preprocessing to mitigate backdoors in the infected model.
 556 Then, during the inference stage, we propose a lightweight
 557 preprocessing function to remove the potential triggers from
 558 the samples. The integration of these techniques can effec-
 559 tively defeat various backdoor threats with different types of
 560 triggers, without any prior adversarial knowledge. Extensive
 561 evaluations show that our method is more robust and com-
 562 prehensive than existing ones.

References

- 563
- 564 Chan, T.-H.; Jia, K.; Gao, S.; Lu, J.; Zeng, Z.; and Ma, Y. 2015. PCANet: A simple deep learning baseline for im-
565 2015. PCANet: A simple deep learning baseline for im-
566 age classification? *IEEE transactions on image processing*
567 24(12): 5017–5032.
- 568 Chen, B.; Carvalho, W.; Baracaldo, N.; Ludwig, H.; Ed-
569 wards, B.; Lee, T.; Molloy, I.; and Srivastava, B. 2018. De-
570 tecting backdoor attacks on deep neural networks by activa-
571 tion clustering. *arXiv preprint arXiv:1811.03728* .
- 572 Chen, H.; Fu, C.; Zhao, J.; and Koushanfar, F. 2019. DeepIn-
573 spect: A Black-box Trojan Detection and Mitigation Frame-
574 work for Deep Neural Networks. In *IJCAI*, 4658–4664.
- 575 Chen, X.; Liu, C.; Li, B.; Lu, K.; and Song, D. 2017. Tar-
576 geted backdoor attacks on deep learning systems using data
577 poisoning. *arXiv preprint arXiv:1712.05526* .
- 578 Chou, E.; Tramèr, F.; Pellegrino, G.; and Boneh, D. 2018.
579 Sentinet: Detecting physical attacks against deep learning
580 systems. *arXiv preprint arXiv:1812.00292* .
- 581 Collobert, R.; and Weston, J. 2008. A unified architecture
582 for natural language processing: Deep neural networks with
583 multitask learning. In *Proceedings of the 25th international*
584 *conference on Machine learning*, 160–167.
- 585 Deng, L.; and Platt, J. C. 2014. Ensemble deep learning for
586 speech recognition. In *Fifteenth Annual Conference of the*
587 *International Speech Communication Association*.
- 588 Du, M.; Jia, R.; and Song, D. 2019. Robust anomaly detec-
589 tion and backdoor attack detection via differential privacy.
590 *arXiv preprint arXiv:1911.07116* .
- 591 Gao, Y.; Xu, C.; Wang, D.; Chen, S.; Ranasinghe, D. C.; and
592 Nepal, S. 2019. Strip: A defence against trojan attacks on
593 deep neural networks. In *Proceedings of the 35th Annual*
594 *Computer Security Applications Conference*, 113–125.
- 595 Gu, T.; Dolan-Gavitt, B.; and Garg, S. 2017. Badnets: Ident-
596 ifying vulnerabilities in the machine learning model supply
597 chain. *arXiv preprint arXiv:1708.06733* .
- 598 Guo, W.; Wang, L.; Xing, X.; Du, M.; and Song, D.
599 2019. Tabor: A highly accurate approach to inspecting and
600 restoring trojan backdoors in ai systems. *arXiv preprint*
601 *arXiv:1908.01763* .
- 602 Han, J.; Shao, L.; Xu, D.; and Shotton, J. 2013. Enhanced
603 computer vision with microsoft kinect sensor: A review.
604 *IEEE transactions on cybernetics* 43(5): 1318–1334.
- 605 He, K.; Zhang, X.; Ren, S.; and Sun, J. 2016. Deep residual
606 learning for image recognition. In *Proceedings of the IEEE*
607 *conference on computer vision and pattern recognition*, 770–
608 778.
- 609 Hosseini, H.; Xiao, B.; and Poovendran, R. 2017. Google’s
610 cloud vision api is not robust to noise. In *2017 16th IEEE*
611 *International Conference on Machine Learning and Appli-*
612 *cations (ICMLA)*, 101–105. IEEE.
- 613 Jin, K.; Zhang, T.; Shen, C.; Chen, Y.; Fan, M.; Lin, C.; and
614 Liu, T. 2020. A unified framework for analyzing and de-
615 tecting malicious examples of dnn models. *arXiv preprint*
616 *arXiv:2006.14871* .
- Kumar, N.; Berg, A. C.; Belhumeur, P. N.; and Nayar, S. K. 617
2009. Attribute and simile classifiers for face verification. 618
In *2009 IEEE 12th international conference on computer vi-* 619
sion, 365–372. IEEE. 620
- Li, S.; Zhao, B. Z. H.; Yu, J.; Xue, M.; Kaafar, D.; and Zhu, 621
H. 2019. Invisible backdoor attacks against deep neural net- 622
works. *arXiv preprint arXiv:1909.02742* . 623
- Li, Y.; Wu, B.; Jiang, Y.; Li, Z.; and Xia, S.-T. 2020a. Back- 624
door Learning: A Survey. *arXiv preprint arXiv:2007.08745* 625
. 626
- Li, Y.; Zhai, T.; Wu, B.; Jiang, Y.; Li, Z.; and Xia, S. 2020b. 627
Rethinking the Trigger of Backdoor Attack. *arXiv preprint* 628
arXiv:2004.04692 . 629
- Liao, C.; Zhong, H.; Squicciarini, A.; Zhu, S.; and Miller, 630
D. 2018. Backdoor embedding in convolutional neural 631
network models via invisible perturbation. *arXiv preprint* 632
arXiv:1808.10307 . 633
- Liberty, E.; Karnin, Z.; Xiang, B.; Rouesnel, L.; Coskun, 634
B.; Nallapati, R.; Delgado, J.; Sadoughi, A.; Astashonok, Y.; 635
Das, P.; et al. 2020. Elastic Machine Learning Algorithms 636
in Amazon SageMaker. In *Proceedings of the 2020 ACM* 637
SIGMOD International Conference on Management of Data, 638
731–737. 639
- Liu, K.; Dolan-Gavitt, B.; and Garg, S. 2018. Fine-pruning: 640
Defending against backdooring attacks on deep neural net- 641
works. In *International Symposium on Research in Attacks,* 642
Intrusions, and Defenses, 273–294. Springer. 643
- Liu, T.; Malcolm, A.; and Xu, J. 2010. Pincushion distor- 644
tion correction in x-ray imaging with an image intensifier. In 645
Fourth International Conference on Experimental Mechan- 646
ics, volume 7522, 75223T. International Society for Optics 647
and Photonics. 648
- Liu, Y.; Lee, W.-C.; Tao, G.; Ma, S.; Aafer, Y.; and Zhang, 649
X. 2019. ABS: Scanning neural networks for back-doors by 650
artificial brain stimulation. In *Proceedings of the 2019 ACM* 651
SIGSAC Conference on Computer and Communications Se- 652
curity, 1265–1282. 653
- Liu, Y.; Ma, S.; Aafer, Y.; Lee, W.-C.; Zhai, J.; Wang, W.; 654
and Zhang, X. 2017. Trojaning attack on neural networks . 655
- Qiao, X.; Yang, Y.; and Li, H. 2019. Defending neural back- 656
doors via generative distribution modeling. In *Advances in* 657
Neural Information Processing Systems, 14004–14013. 658
- Stallkamp, J.; Schlipsing, M.; Salmen, J.; and Igel, C. 2012. 659
Man vs. computer: Benchmarking machine learning algo- 660
rithms for traffic sign recognition. *Neural networks* 32: 323– 661
332. 662
- Tran, B.; Li, J.; and Madry, A. 2018. Spectral signatures in 663
backdoor attacks. In *Advances in Neural Information Pro-* 664
cessing Systems, 8000–8010. 665
- Wang, B.; Yao, Y.; Shan, S.; Li, H.; Viswanath, B.; Zheng, 666
H.; and Zhao, B. Y. 2019. Neural cleanse: Identifying and 667
mitigating backdoor attacks in neural networks. In *2019* 668
IEEE Symposium on Security and Privacy (SP), 707–723. 669
IEEE. 670

671 Xu, X.; Wang, Q.; Li, H.; Borisov, N.; Gunter, C. A.; and Li,
672 B. 2019. Detecting AI Trojans Using Meta Neural Analysis.
673 *arXiv preprint arXiv:1910.03137* .

674 Yao, Y.; Li, H.; Zheng, H.; and Zhao, B. Y. 2019. Latent
675 backdoor attacks on deep neural networks. In *Proceedings of*
676 *the 2019 ACM SIGSAC Conference on Computer and Com-*
677 *munications Security*, 2041–2055.

678 Zeiler, M. D. 2012. Adadelta: an adaptive learning rate
679 method. *arXiv preprint arXiv:1212.5701* .

680 Zeng, Y.; Qiu, H.; Memmi, G.; and Qiu, M. 2020. A
681 Data Augmentation-based Defense Method Against Ad-
682 versarial Attacks in Neural Networks. *arXiv preprint*
683 *arXiv:2007.15290* .