

---

# Fighting Gradients with Gradients: Dynamic Defenses against Adversarial Attacks

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Adversarial attacks optimize against models to defeat defenses. Existing defenses  
2 are static, and stay the same once trained, even while attacks change. We argue  
3 that models should fight back, and optimize their defenses against attacks at test  
4 time. We propose dynamic defenses, to adapt the model and input during testing,  
5 by defensive entropy minimization (dent). Dent alters testing, but not training, for  
6 compatibility with existing models and train-time defenses. Dent improves the  
7 robustness of adversarially-trained defenses and nominally-trained models against  
8 white-box, black-box, and adaptive attacks on CIFAR-10/100 and ImageNet. In  
9 particular, dent boosts state-of-the-art defenses by 20+ points absolute against  
10 AutoAttack on CIFAR-10 at  $\epsilon_\infty = 8/255$ .

## 11 1 Introduction: Attack, Defend, and Then?

12 Deep networks are vulnerable to adversarial attacks: input perturbations that alter natural data to  
13 cause errors or exploit predictions [54]. As deep networks are deployed in real systems, these attacks  
14 are real threats [63], and so defenses are needed. The challenge is that every new defense is followed  
15 by a new attack, in a loop [56]. The strongest attacks, armed with gradient optimization, update to  
16 circumvent defenses that do not. Such iterative attacks form an even tighter loop to ensnare defenses.  
17 In a cat and mouse game, the mouse must keep moving to survive.

18 Current defenses, deterministic or stochastic, stand still: once trained, they are *static* and do not adapt  
19 during testing. Adversarial training [18, 30] learns from attacks during training, but cannot learn  
20 from test data. Stochastic defenses alter the network [11] or input [20, 7], but their randomness is  
21 independent of test data. Static defenses do not adapt, and so they may fail as attacks update.

22 Our *dynamic* defense fights adversarial updates with defensive updates by adapting during testing  
23 (Figure 1). In fact, our defense updates on every input, whether natural or adversarial. Our defense  
24 objective is entropy minimization, to maximize model confidence, so we call our method *dent* for  
25 defensive entropy. Our updates rely on gradients and batch statistics, inspired by test-time adaptation  
26 approaches [53, 43, 28, 29, 58]. In pivoting from training to testing, dent is able to keep changing, so  
27 the attacker never hits the same defense twice. Dent has the last move advantage, as its update always  
28 follows each attack.

29 Dent connects adversarial defense and domain adaptation, which share an interest in the sensitivity of  
30 deep networks to input shifts. Just as models fail on adversarial attacks, they fail on natural shifts  
31 like corruptions. Adversarial data is a particularly hard shift, as evidenced by the need for more  
32 parameters and optimization for adversarial training [30], and its negative side effect of reducing

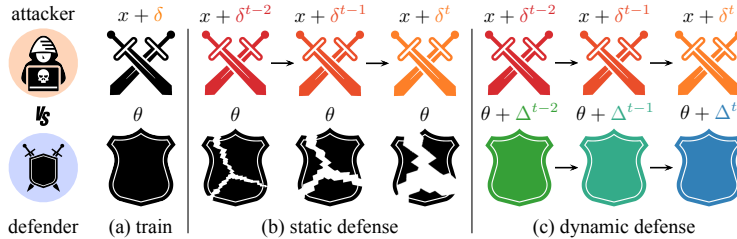


Figure 1: Attacks optimize the input  $x + \delta$  against the model  $\theta$ . Adversarial training optimizes  $\theta$  for defense (a), but attacks update during testing while  $\theta$  does not (b). Our *dynamic* defense improves robustness by adapting  $\theta + \Delta$  during testing (c), so the attack cannot hit the same defense twice.

33 accuracy on natural data [52, 65]. Faced with these difficulties, we turn to adaptation, and change our  
 34 focus to testing, rather than training more still.

35 Experiments evaluate dent against white-box attacks (APGD, FAB), black-box attack (Square), and  
 36 adaptive attacks that are aware of its updates. Dent boosts state-of-the-art adversarial training defenses  
 37 on CIFAR-10 by 20+ points against AutoAttack [9] at  $\epsilon_\infty = 8/255$ . Ablations inspect the effects of  
 38 iteration, parameterization, and batch size. Our code is included in the supplement.

### 39 Our contributions

- 40 • We highlight an opportunity for dynamic defense: the last move advantage.
- 41 • We propose the first fully test-time dynamic defense: dent adapts both the model and input  
 42 during testing without needing to alter training.
- 43 • Dent augments state-of-the-art adversarial training methods, improving robustness by 30%  
 44 relative, and tops the AutoAttack leaderboard by 15+ points.
- 45 • We devise two adaptive attacks against dent: denying updates and mixing batches.

## 46 2 Related Work

47 **Adversarial Defense** For adaptive adversaries, which change in response to defenses, it is natural  
 48 to consider dynamic defenses, which adapt in turn. Evans et al. [14] explain dynamic defenses are  
 49 promising in principle but caution they may not be effective in practice. Their analysis concerns  
 50 randomized defenses, which do change, but their randomization does not adapt to the input. We argue  
 51 for dynamic defenses that depend on the input to keep adapting along with the attacks. Goodfellow  
 52 [17] supports dynamic defenses for similar reasons, but does not develop a specific defense. We  
 53 demonstrate the first defense to optimize the model and input during testing for improved robustness.

54 Most defenses for deep learning focus on first-order adversaries [18, 30] that are equipped with  
 55 gradient optimization but constrained by  $\ell_p$ -norm bounds. Adversarial training and randomization  
 56 are the most effective defenses against such attacks, but are nevertheless limited, as they are fixed  
 57 during testing. Adversarial training [18, 30] trains on attacks, but a different or stronger adversary  
 58 (by norm or bound) can overcome the trained defense [46, 55]. Randomizing the input [37, 7, 32]  
 59 or network [11] requires the adversary to optimize in expectation [3], but can still fail with more  
 60 iterations. Furthermore, these defenses gain adversarial robustness by sacrificing accuracy on natural  
 61 data. Dent adapts during testing to defend against various attacks without more harm to natural  
 62 accuracy.

63 Generative, self-supervised, and certified defenses try to align testing with training but are still  
 64 static. Generative defenses optimize the input w.r.t. autoregressive [50], GAN [42], or energy [23]  
 65 models, but the models do not adapt, and may be attacked by approximating their gradients [3].  
 66 Self-supervised defenses optimize the input w.r.t. auxiliary tasks [49], but again the models do not  
 67 adapt. Certified defenses [7, 66] guarantee robustness within their training scope, but are limited  
 68 to small perturbations by specific types of attacker during testing. Changing data distributions or

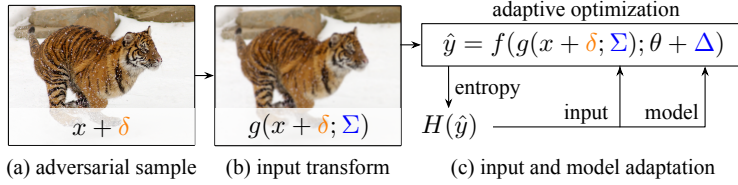


Figure 2: Dent adapts the model and input to minimize the entropy of the prediction  $H(\hat{y})$ . The model  $f$  is adapted by a constrained update  $\Delta$  to the parameters  $\theta$ . The input is adapted by smoothing  $g$  with parameters  $\Sigma$ . Dent updates batch-by-batch during testing.

69 adversaries requires re-training all of these defenses. Dent adapts during testing, without requiring  
 70 (re-)training, and is the only method to update the model itself against attack.

71 **Domain Adaptation** Domain adaptation mitigates input shifts between the source (train) and target  
 72 (test) to maintain model accuracy [34, 41]. Adversarial attacks are such a shift, and adversarial  
 73 error is related to natural generalization error [51, 15]. How then can adaptation inform dynamic  
 74 defense? Train-time adaptation is static, like adversarial training, with the same issues of capacity,  
 75 optimization, and re-computation when the data/adversary changes. We instead turn to test-time  
 76 adaptation methods.

77 Test-time adaptation keeps updating the model as the data changes. Model parameters and statistics  
 78 can be updated by self-supervision [53], normalization [43], and entropy minimization [58]. These  
 79 methods improve robustness to natural corruptions [22], but their effect on adversarial perturbations  
 80 is not known. We base our defense on entropy minimization as it enables optimization during testing  
 81 without altering model architecture or training (as needed for self-supervision). For defense, we (1)  
 82 extend the parameterization of adaptation with model and input transformations, (2) optimize for  
 83 additional iterations, and (3) investigate usage on data that is adversarial, natural, or mixed. We are  
 84 the first to report test-time model adaptation improves robustness to adversarial perturbations.

85 **Dynamic Inference** A *dynamic* model conditionally changes inference for each input, while a  
 86 *static* model unconditionally fixes inference for all inputs. There are various dynamic inference  
 87 techniques, with equally varied goals, such as expressivity with more parameters or efficiency with  
 88 less computation. All static models are alike; each dynamic model is dynamic in its own way.

89 Selection techniques learn to choose a subset of components [1, 57]. Halting techniques learn to  
 90 continue or end computation [19, 59]. Mixing techniques learn to combine parameters [47, 33, 62].  
 91 Implicit techniques learn to iteratively update [6, 4]. While these methods learn to adapt during  
 92 *training*, our method keeps adapting by directly optimizing during *testing*.

### 93 3 Dynamic Defense by Test-Time Adaptation

94 Adversarial attacks optimize against defenses at test time, so defenses should fight back, and counter-  
 95 optimize against attacks. Defensive entropy minimization (dent) does exactly this for dynamic  
 96 defense by test-time adaptation.

97 In contrast to many existing defenses, dent alters testing, but not training. Dent only needs differenti-  
 98 able parameters for gradient optimization and probabilistic predictions for entropy measurement.  
 99 As such, it applies to both adversarially-trained and nominally-trained models.

#### 100 3.1 Preliminaries on Attacks and Defenses

101 Let  $x \in \mathbb{R}^d$  and  $y \in \{1, \dots, C\}$  be an input sample and its corresponding ground truth. Given a  
 102 model  $f(\cdot; \theta): \mathbb{R}^d \rightarrow \mathbb{R}^C$  parameterized by  $\theta$ , the goal of the adversary is to craft a perturbation  
 103  $\delta \in \mathbb{R}^d$  such that the perturbed input  $\tilde{x} = x + \delta$  causes a prediction error  $f(x + \delta; \theta) \neq y$ .

104 A targeted attack aims for a specific prediction of  $y'$ , while an untargeted attack seeks any incorrect  
 105 prediction. The perturbation  $\delta$  is constrained by a choice of  $\ell_p$  norm and threshold  $\epsilon$ :  $\{\delta \in \mathbb{R}^d \mid$   
 106  $\|\delta\|_p < \epsilon\}$ . We consider the two most popular norms for adversarial attacks:  $\ell_\infty$  and  $\ell_2$ .

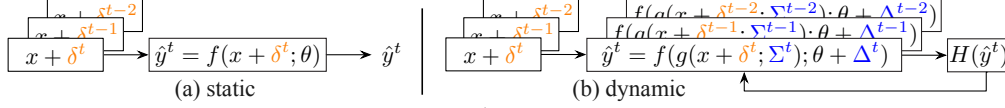


Figure 3: The adversary optimizes its attacks  $\delta^{1 \dots t}$  against the model  $f$ . Static defenses (left) do not adapt, and are vulnerable to persistent, iterative attacks. Our dynamic defenses (right) do adapt, and update their parameters  $\Delta, \Sigma$  each time the adversary updates its attack  $\delta$ .

107 Adversarial training is a standard defense, formulated by Madry et al. [30] as a saddle point problem,

$$\operatorname{argmin}_{\theta} \mathbb{E}_{(x,y)} \max_{\delta} L(f(x + \delta; \theta), y), \quad (1)$$

108 which the model minimizes and the adversary maximizes with respect to the loss  $L(\hat{y}, y)$ , such as  
 109 cross-entropy for classification. The adversary iteratively optimizes  $\delta$  by projected gradient descent  
 110 (PGD), a standard algorithm for constrained optimization, for each step  $t$  via

$$\delta^t = \Pi_p(\delta^{t-1} + \alpha \cdot \operatorname{sign}(\nabla_{\delta^{t-1}} L(f(x + \delta^{t-1}; \theta), y))), \quad (2)$$

111 for projection  $\Pi_p$  onto the norm ball for  $\ell_p < \epsilon$ , step size hyperparameter  $\alpha$ , and random initialization  
 112  $\delta^0$ . The model optimizes  $\theta$  against  $\delta$  to minimize the loss of its predictions on perturbed inputs. This  
 113 is accomplished by augmenting the training set with adversarial inputs from PGD attack.

114 Adversarial training is state-of-the-art, but static. Dynamic defenses offer to augment its robustness.

### 115 3.2 Defensive Entropy Minimization

116 Defensive entropy minimization (dent) counters attack updates with defense updates. While adver-  
 117 saries optimize to cross decision boundaries, entropy minimization optimizes to distance predictions  
 118 from decision boundaries, interfering with attacks. As the adversary optimizes its perturbation  $\delta$ , dent  
 119 optimizes its adaptation  $\Delta, \Sigma$ . Figure 2 shows dent’s model ( $\Delta$ ) and input ( $\Sigma$ ) updates.

120 Dent is dynamic because both  $\Delta, \Sigma$  depend on the testing data, whether natural  $x$  or adversarial  
 121  $x + \delta$ . On the contrary, static defenses depend only on training data through the model parameters  $\theta$ .  
 122 Figure 3 contrasts static and dynamic defenses across the steps of attack optimization.

123 **Entropy Objective** Test-time optimization requires an unsupervised objective. Following tent [58],  
 124 we adopt entropy minimization as our adaptation objective. Specifically, our defense objective is to  
 125 minimize the Shannon entropy [45]  $H(\hat{y})$  of the model prediction during testing  $\hat{y} = f(x; \theta)$  for the  
 126 probability  $\hat{y}_c$  of class  $c$ :

$$H(\hat{y}) = - \sum_{c \in \{1, \dots, C\}} p(\hat{y}_c) \log p(\hat{y}_c) \quad (3)$$

127 **Adaptation Parameters** Dent adapts the model by  $\Delta$  and input by  $\Sigma$  (Figure 2). For the model, dent  
 128 adapts affine scale  $\gamma$  and shift  $\beta$  parameters by gradient updates and adapts mean  $\mu$  and variance  $\sigma^2$   
 129 statistics by estimation. These are a small portion of the full model parameters  $\theta$ , in only the batch  
 130 normalization layers [25]. However, they are effective for conditioning a model on changes in the  
 131 task [33] or data [43, 58]. For the input, dent updates Gaussian smoothing  $g$  by gradient updates  
 132 of the parameter  $\Sigma$ , while adjusting the filter size for efficiency [48]. This controls the degree of  
 133 smoothing dynamically, unlike defense by static smoothing [7].

134 In standard models the scale  $\gamma$  and shift  $\beta$  parameters are shared across inputs, and so adaptation  
 135 updates batch-wise. For further adaptation, dent can update sample-wise, with different affine  
 136 parameters for each input. In this way, it adapts more than prior test-time adaptation methods with  
 137 batch-wise parameters [58, 43].

138 Our model and input parameters are differentiable, so end-to-end optimization coordinates them  
 139 against attacks as layered defenses. This coordination is inspired by CyCADA [24], for domain  
 140 adaptation, but dent differs in its purpose and its unified loss. CyCADA also optimizes input and  
 141 model transformations but does so in parallel with separate losses. Our defensive optimization is  
 142 joint and shares the same loss.

Table 1: Dent boosts the robustness of adversarial training on CIFAR-10 against AutoAttack. Adversarial training is static, but dent is dynamic, and adapts during testing. Dent adapts batch-wise, while dent+ adapts sample-wise, surpassing the state-of-the-art for static defense at robustbench.github.io.

ACCURACY(%)	NATURAL	ADVERSARIAL		
		STATIC	DENT	DENT+
$\epsilon_\infty = 8/255$				
CARMON ET AL. [5]	89.6	59.5	74.7	<b>82.3</b>
SEHWAG ET AL. [44]	84.4	54.4	61.2	<b>75.2</b>
WONG ET AL. [60]	83.3	43.2	52.3	<b>71.8</b>
DING ET AL. [12]	88.0	41.4	47.6	<b>64.4</b>
$\epsilon_2 = 0.5$				
SEHWAG ET AL. [44]	89.5	73.4	77.8	<b>85.7</b>
RICE ET AL. [38]	88.7	67.7	69.7	<b>81.3</b>
RONY ET AL. [39]	89.1	66.4	73.4	<b>85.3</b>
DING ET AL. [12]	88.0	66.1	70.3	<b>82.8</b>

143 **Update Algorithm** In summary, when the adversary attacks with perturbation  $\delta^t$ , our dynamic  
 144 defense reacts with  $\Sigma^t, \Delta^t$ . The parameters of the model  $f$  and smoothing  $g$  are updated by  
 145  $\operatorname{argmin}_{\Sigma, \Delta} H(f(g(x + \delta; \Sigma); \theta + \Delta))$  through test-time optimization. At each step, dent estimates  
 146 the normalization statistics  $\mu, \sigma$  and then updates the parameters  $\gamma, \beta, \Sigma$  by the gradient of entropy  
 147 minimization. Figure 3 contrasts static defenses and dynamic defenses that update like dent.

148 Dent adapts on batches rather than samples. Batch-wise adaptation stabilizes optimization for entropy  
 149 minimization. The defense parameters reset between batches.

150 **Discussion** The purpose of a dynamic defense is to move when the adversary moves. When the  
 151 adversary submits an attack  $x + \delta^t$ , the defense counters with  $\Delta^t$ . In this way, the defense has the  
 152 last move, and therefore an advantage.

153 Our dynamic defense changes the model, and therefore its gradients, but differs from gradient  
 154 obfuscation [3]. Our defense does not rely on (1) shattered gradients, as the update does not cause  
 155 non-differentiability or numerical instability; (2) stochastic gradients, as the update is deterministic  
 156 given the input, model, and prior updates; nor (3) exploding/vanishing gradients, as the update  
 157 improves robustness with even a single step (although more steps are empirically better).

158 Dent forces the attack to rely on a *stale* gradient, as  $\delta^t$  follows  $\Delta^{t-1}$ , while the model adapts by  $\Delta^t$ .

## 159 4 Experiments

160 We evaluate dent against white-box, black-box, and adaptive attacks with a variety of static defenses  
 161 and datasets. For attacks, we choose the AutoAttack [9] benchmark, which includes four attack types  
 162 spanning white-box/gradient and black-box/query attacks. For static defenses, we choose strong and  
 163 recent adversarial training methods, and we also experiment with nominally trained models. For  
 164 datasets, we evaluate dent on CIFAR-10/CIFAR-100 [27], as they are popular datasets for adversarial  
 165 robustness, and ImageNet [40], as it is a large-scale dataset.

166 We ablate the choice of model/input adaptation, parameterization, and the number of updates.

### 167 4.1 Setup

168 **Metrics** We score natural accuracy on the regular test data  $x$  and adversarial accuracy on the perturbed  
 169 test data  $x + \delta$ . Each is measured as percentage accuracy (higher is better). We report the worst-case  
 170 adversarial accuracy across attacks.

171 **Test-time Optimization** We optimize batch-wise  $\Delta$  (dent) and sample-wise  $\Delta$  (dent+). Dent updates  
 172 by Adam [26] with learning rate 0.001. Dent+ updates by AdaMod [13] with learning rate 0.006.  $\Sigma$

173 updates use learning rate 0.25. All updates use batch size 128 and no weight decay. Dent+ regularizes  
174 updates by information maximization [16, 29]. We tuned update hyperparameters against PGD  
175 attacks. Please see the code for exact settings.

176 **Architecture** For comparison with existing defenses, we keep the architecture and training the same,  
177 and simply load the public reference models provided by RobustBench [10]. For analysis and ablation  
178 experiments, we define a residual net with 26 layers and a width multiplier of 4 (ResNet-26-4) [21, 64],  
179 following prior work on adaptation [53, 58].

## 180 4.2 Attack Types & Threat Model

181 We evaluate standard white-box and black-box attacks with adversarially-trained models (Section 4.3)  
182 and nominally-trained models (Section 4.4), as well as dent-specific adaptive attacks (Section 4.5).

183 We primarily evaluate against AutoAttack’s ensemble of:

- 184 1. APGD-CE [30, 9], an untargeted white-box attack by cross-entropy,
- 185 2. APGD-DLR [9], a targeted white-box attack with a shift and scale invariant loss,
- 186 3. FAB [8], a targeted white-box attack for minimum-norm perturbation,
- 187 4. Square Attack [2], an untargeted black-box attack with square-shaped updates.

188 These attacks are cumulative, so a defense is only successful if it holds against each type. Following  
189 convention, we evaluate  $\ell_\infty$  attacks with  $\epsilon_\infty = 8/255$  and  $\ell_2$  attacks with  $\epsilon_2 = 0.5$ . This is the  
190 standard evaluation adopted by the popular RobustBench benchmark [10].

191 We devise and experiment with two adaptive attacks against dent and its dynamic updates. The first  
192 interferes with adaptation by denying updates: it optimizes offline against  $\theta$  without  $\Delta, \Sigma$  updates.  
193 The second interferes with adaptation by mixing data: it combines adversarial data and natural data  
194 in the same batch. Both are specific to dent to complement our general evaluation by AutoAttack.

195 These attacks fall under the usual white-box threat model. The adversary has full access to the  
196 classifier, including its architecture and parameters, and the defense, such as dent’s adaptation  
197 parameters and statistics. With this access the adversary chooses an attack for each input, but it  
198 cannot choose the inputs (the test set is fixed).

199 We include one additional requirement: dent assumes access to test *batches* rather than individual  
200 test *samples*. While independent, sample-wise defense is ideal for simplicity and latency, batch  
201 processing is not impractical. For example, cloud deployments of deep learning batch inputs for  
202 throughput efficiency, and large-scale systems handle many inputs per unit time [31].

203 The supplementary material covers more attacks, including AutoAttack Plus and Boundary, to confirm  
204 that AutoAttack is a sufficient measure of robustness.

## 205 4.3 Dynamic Defense of Adversarial Training

206 We extend static adversarial training defenses with dynamic updates by dent. Compared to nominal  
207 training, adversarial training achieves higher adversarial accuracy but lower natural accuracy. The  
208 purpose of dent is to improve adversarial accuracy without further harming natural accuracy.

209 **Dent improves state-of-the-art defenses.** Table 1 shows state-of-the-art adversarial training defenses  
210 [5, 44, 39, 38, 60, 12] with and without dynamic defense by dent. Note that dent does not specialize  
211 to the choice of norm or bound, unlike adversarial training, but instead adapts to each attack during  
212 testing. In each case, dent significantly improves adversarial accuracy and maintains natural accuracy.

213 Dent updates batch-wise for 30 steps. Dent+ is more robust in fewer steps by sample-wise adaptation.  
214 With sample-wise  $(\gamma, \beta)$  parameters, dent+ needs only six steps to reach an adversarial accuracy  
215 within 90% of the natural accuracy. These experiments only include model adaptation of  $\Delta$ , without  
216 input adaptation of  $\Sigma$ , as we found it unnecessary when combined with adversarial training.

Table 2: AutoAttack includes four attack types, and dent improves robustness to each on CIFAR-10 against  $\ell_\infty$  attacks. We evaluate without dent (-) and with dent (+).

ACCURACY(%)	APGD-CE		APGD-DLR		FAB		SQUARE	
	-	+	-	+	-	+	-	+
WONG ET AL. [60]	45.9	57.6	43.2	52.3	43.2	52.3	43.2	52.3
DING ET AL. [12]	50.1	60.2	41.6	48.0	41.5	47.7	41.4	47.6

Table 3: Ablation of model adaptation ( $\Delta$ ), input adaptation ( $\Sigma$ ), and steps on the accuracy of a nominally-trained model with dent.

$\Delta$	$\Sigma$	STEP	TIME	NATURAL	ADVERSARIAL	
					$\epsilon_\infty = \frac{1.5}{255}$	$\epsilon_2 = 0.2$
×	NONE	0	1.0×	95.6	8.8	9.2
✓	NONE	1	3.6×	95.6	15.0	13.5
×	STAT.	0	1.0×	86.2	25.8	23.6
✓	STAT.	1	3.6×	86.3	27.5	24.4
✓	STAT.	10	25.9×	86.3	37.6	30.9
✓	DYNA.	10	26.1×	92.5	45.4	36.5

217 **Dent helps across attack types.** Table 2 evaluates dent against each attack in the AutoAttack  
 218 ensemble. Dent improves robustness to each attack type. We report the worst case across these types  
 219 in the remainder of our experiments.

220 **Dent helps across datasets and architectures.** We experiment on ImageNet to check scalability.  
 221 We evaluate the defense of Wong et al. [60], one of few defenses that scales to this dataset, against  
 222 strong  $\ell_\infty$ -PGD attacks with 30 iterations, step size of 0.1, and five random starts. Dent improves the  
 223 adversarial accuracy by 14 points against PGD at  $\epsilon_\infty = 4/255$  and natural accuracy by 23 points.

224 Table ?? in the supplement confirms improvement across more defenses, architectures, and datasets.

#### 225 4.4 Dynamic Defense of Nominal Training

226 Dent improves the adversarial accuracy of off-the-shelf, nominally-trained models. As dent does not  
 227 assume adversarial training, it can apply to various models at test time.

228 For nominal training, we exactly follow the CIFAR reference training in pycls [35, 36] with ResNet-  
 229 26-4/ResNet-32-10 architectures. Briefly, we train by stochastic gradient descent (SGD) for 200  
 230 epochs with batch size 128, learning rate 0.1 and decay 0.0005, momentum 0.9, and a half-period  
 231 cosine schedule.

232 We evaluate against  $\ell_\infty$  and  $\ell_2$  AutoAttack attacks on CIFAR-10. As the nominally-trained models  
 233 have no static defense, we constrain the adversaries to smaller  $\epsilon$  perturbations.

234 **Dent defends nominally-trained models.** Table 3 inspects how each part of dent affects adversarial  
 235 accuracy and natural accuracy. When applying dent to nominally-trained models, model adaptation  
 236 through  $\Delta$  is further helped by input adaptation through  $\Sigma$ . In just a single step, the  $\Delta$  update  
 237 improves adversarial accuracy without affecting natural accuracy. from 8.8% to 15.0% against  $\ell_\infty$   
 238 attacks with just a single step. With 10 steps, and  $\Sigma$  adaptation, dent improves the model’s adversarial  
 239 accuracy to 45.4% against  $\ell_\infty$  attacks and 36.5% against  $\ell_2$  attacks. In total, dent boosts  $\ell_\infty$  and  $\ell_2$   
 240 adversarial accuracy by almost 40 and 30 points while only sacrificing 3 points of natural accuracy.  
 241 Dent delivers this boost at test-time, without re-training.

242 **Input adaptation helps preserve natural accuracy.** Gaussian smoothing significantly improves  
 243 adversarial accuracy. This agrees with prior work on denoising by optimization [20] or randomized  
 244 smoothing [7]. Tuned as a fixed hyperparameter, smoothing helps adversarial accuracy but hurts  
 245 natural accuracy. Optimized end-to-end, our dynamic smoothing reduces the natural accuracy gap.  
 246 On natural data, the learned  $\Sigma$  for the blur decreases to approximate the identity transformation.

Table 4: Adaptive attack by denying updates. We transfer attacks from static models to dent and then evaluate nominal and adversarial training [30] against  $\ell_\infty$  and  $\ell_2$  AutoAttack. Attacks break the static models (static-static), but fail to transfer to our dynamic defense (static-dent).

	NOMINAL		ADVERSARIAL	
	$\epsilon_\infty = \frac{1.5}{255}$	$\epsilon_2 = 0.2$	$\epsilon_\infty = \frac{8}{255}$	$\epsilon_2 = 0.5$
STATIC-STATIC	11.6	11.0	42.0	44.1
STATIC-DENT	82.5	81.6	50.0	50.2

Table 5: Adaptive attack by mixing adversarial and natural data. We report the adversarial accuracy on mixed batches, from low to high amounts of adversarial data. Dent improves on adversarial training (43.8%) across mixing proportions within 10 steps.

$\mu, \sigma$	STEP	1	10%	25%	50%	75%	90%
$\times$	1	-	43.4	43.2	44.0	44.2	43.8
$\times$	10	62.4	51.2	49.6	48.7	48.7	47.6
$\checkmark$	1	-	41.7	41.4	43.2	44.1	44.7
$\checkmark$	10	54.9	47.6	47.7	49.7	50.6	50.9

247 **4.5 Adaptive Attacks on Dent Updates**

248 We adaptively attack dent through its use of adaptation by (1) denying updates and (2) mixing batches.  
 249 To deny updates, we attack the static model offline by optimizing against  $\theta$  without  $\Delta, \Sigma$  updates,  
 250 then submit this attack to dent. This attempts to short circuit adaptation by disrupting the first update  
 251 with a sufficiently strong perturbation. To mix batches, we mix adversarial and natural data in the  
 252 same batch. This attempts to prevent adaptation by aligning batch statistics with natural data.

253 **Denying Updates** The aim of this attack is to defeat adaptation on the first move, before dent can  
 254 update to counter it. We optimize against the static model alone to prevent defensive optimization  
 255 until adversarial optimization is complete. Under this attack, the input to dent is the final perturbation  
 256 derived by adversarial attack against the static model.

257 We examine whether these offline perturbations can disrupt adaptation. Table 4 shows that dent can  
 258 still defend against this attack. This suggests that updating, and having the last move, remains an  
 259 advantage for our dynamic defense.

260 **Mixing Batches** Dent adapts batch-wise, with the underlying assumption that one shared transfor-  
 261 mation can defend the whole batch. We challenge this assumption by evaluating mixed batches of  
 262 adversarial and natural data. In Table 5, we vary the ratio of adversarial and natural data in each batch  
 263 and measure accuracy on the adversarial portion.

264 At the extreme, we consider an adaptive attack with only one adversarial input per batch. Specifically,  
 265 we batch one adversarial input with 15 natural inputs randomly chosen from the test set. This adaptive  
 266 attack aims to reduce adaptation by the dynamic defense, as natural inputs do not need adaptation.

267 Dent is generally robust to batch mixing, and improves over adversarial training in 10 steps or less.

268 **4.6 Ablations & Analysis**

269 **More updates deliver more defense.** The number of steps can balance defense and computation.  
 270 Table 6 shows that more steps offer stronger defense for both dent and dent+. However, more steps  
 271 do nevertheless require more computation: ten-step optimization takes  $25.9\times$  more operations than  
 272 the static model (Table 3). As a plus, dent+ is not only more robust, but also more efficient in needing  
 273 fewer steps. Note that the computational difference between dent and dent+ is negligible, as the  
 274 adaptation parameters are such a small fraction of the model.

275 **Model adaptation updates depend on the attack type.** Dent adapts by adjusting normalization  
 276 statistics and affine transformation parameters. Dent can fix or update the normalization statistics  
 277 ( $\mu, \sigma$ ) by using static training statistics ( $\times$ ) or dynamic testing statistics ( $\checkmark$ ); Dent can fix or update  
 278 the affine parameters ( $\gamma, \beta$ ) by not taking gradients ( $\times$ ) or applying gradient updates ( $\checkmark$ ). Table 7  
 279 compares each combination: affine updates always help, but both updates together hurt  $\ell_2$  robustness.

280 **Batch size** We analyze dent’s sensitivity to batch size and focus on small batch sizes. Some real-world  
 281 tasks, such as autonomous driving, naturally provide a small batch of inputs (from consecutive video



Table 6: Dynamic defenses can trade computation and adaptation. More steps are more robust on CIFAR-10 with  $\ell_\infty$  AutoAttack. Dent+ reaches higher adversarial accuracy in fewer steps.

DENT	STEPS			
	0	20	30	40
CARMON ET AL. [5]	59.5	68.3	74.7	76.1
WONG ET AL. [60]	43.2	48.2	52.3	55.1
DING ET AL. [12]	41.4	45.4	47.6	48.7
DENT+	0	1	3	6
DING ET AL. [12]	41.4	46.5	57.7	64.4

Table 7: Ablation of model adaptation with and without normalization statistics ( $\mu, \sigma$ ) and affine parameters ( $\gamma, \beta$ ) updates.

ACCURACY(%)		NOMINAL		ADVERSARIAL	
$\mu, \sigma$	$\gamma, \beta$	$\epsilon_\infty = \frac{1.5}{255}$	$\epsilon_2 = 0.2$	$\epsilon_\infty = \frac{8}{255}$	$\epsilon_2 = 0.5$
×	×	8.8	9.2	43.8	47.3
✓	×	11.7	11.2	41.8	44.1
×	✓	16.8	16.2	49.9	57.3
✓	✓	21.2	15.2	50.4	53.0

Table 8: Sensitivity analysis of batch size and adversarial accuracy with dent. With static batch statistics (×), small batch sizes are better. With dynamic batch statistics (✓), small batch sizes are worse.

$\mu, \sigma$	TYPE	1	2	4	8	16	32	64
×	NAT.	85.9	86.0	85.9	85.9	86.1	86.1	86.2
×	ADV.	70.4	69.5	67.8	65.3	61.9	58.6	55.1
✓	NAT.	11.1	68.1	76.3	80.9	83.4	84.9	85.8
✓	ADV.	5.8	35.9	48.3	53.0	55.3	54.4	52.9

282 frames or various cameras, for example), and so we confirm that dent can maintain robustness on  
 283 such small batches. Table 8 varies batch sizes to check dent’s natural and adversarial accuracy.

## 284 5 Discussion

285 In advocating for dynamic defenses, we hope that test-time updates can help level the field for attacks  
 286 and defenses. Our proposed defensive entropy method takes a first step by countering adversarial  
 287 optimization with defensive optimization over the model and input. While more test-time computation  
 288 is needed for the back-and-forth iteration of attacks and defenses, the cost of defense scales with the  
 289 cost of attack, and some use cases may prefer slow and strong to fast and wrong.

290 **Limitations** Dent depends on batches to adapt, especially for fully test-time defense without ad-  
 291 versarial training. It also relies on a particular choice of model and input parameters. A different  
 292 objective could possibly lessen its dependence on batch size and reliance on constrained updates.  
 293 More generally, dynamic defenses may present difficulties for certification or deployment, as they  
 294 could drift. Along with how to update, improved defenses could investigate when to reset, or how to  
 295 batch inputs for joint optimization.

296 **Benchmarking** Standardized benchmarking, by AutoAttack and RobustBench for example, drives  
 297 progress by competition and empirical corroboration. Dent brings adversarial accuracy on their  
 298 benchmark within 90% of natural accuracy for three of the most accurate methods tested [5, 61, 12].  
 299 This is encouraging, but more research is needed to fully characterize dynamic defenses like dent.  
 300 However, RobustBench is designed for static defenses, and disqualifies dent by its rule against  
 301 test-time optimization. Continued progress could depend on a new benchmark to standardize rules  
 302 for how attacks and defenses alike may adapt.

303 By fighting gradients with gradients, dent shows the potential for dynamic defenses to update and  
 304 counter adversarial attacks. The next steps—by attacks and defenses—will tell.

## References

- 305 [1] Andreas, J., Rohrbach, M., Darrell, T., and Klein, D. Neural module networks. In *CVPR*, 2016.
- 306 [2] Andriushchenko, M., Croce, F., Flammarion, N., and Hein, M. Square attack: a query-efficient black-box  
307 adversarial attack via random search. In *ECCV*, 2020.
- 308 [3] Athalye, A., Carlini, N., and Wagner, D. Obfuscated gradients give a false sense of security: Circumventing  
309 defenses to adversarial examples. In *ICML*, 2018.
- 310 [4] Bai, S., Kolter, J. Z., and Koltun, V. Deep equilibrium models. In *NeurIPS*, 2020.
- 311 [5] Carmon, Y., Ragunathan, A., Schmidt, L., Duchi, J. C., and Liang, P. S. Unlabeled data improves  
312 adversarial robustness. In *NeurIPS*, 2019.
- 313 [6] Chen, R. T., Rubanova, Y., Bettencourt, J., and Duvenaud, D. Neural ordinary differential equations. In  
314 *NeurIPS*, 2018.
- 315 [7] Cohen, J., Rosenfeld, E., and Kolter, Z. Certified adversarial robustness via randomized smoothing. In  
316 *ICML*, 2019.
- 317 [8] Croce, F. and Hein, M. Minimally distorted adversarial examples with a fast adaptive boundary attack. In  
318 *ICML*, 2020.
- 319 [9] Croce, F. and Hein, M. Reliable evaluation of adversarial robustness with an ensemble of diverse parameter-  
320 free attacks. In *ICML*, 2020.
- 321 [10] Croce, F., Andriushchenko, M., Sehwag, V., Flammarion, N., Chiang, M., Mittal, P., and Hein, M.  
322 Robustbench: a standardized adversarial robustness benchmark. *arXiv preprint arXiv:2010.09670*, 2020.
- 323 [11] Dhillon, G. S., Azizzadenesheli, K., Bernstein, J. D., Kossaifi, J., Khanna, A., Lipton, Z. C., and Anandku-  
324 mar, A. Stochastic activation pruning for robust adversarial defense. In *ICLR*, 2018.
- 325 [12] Ding, G. W., Sharma, Y., Lui, K. Y. C., and Huang, R. Mma training: Direct input space margin  
326 maximization through adversarial training. In *ICLR*, 2020.
- 327 [13] Ding, J., Ren, X., Luo, R., and Sun, X. An adaptive and momental bound method for stochastic learning.  
328 *arXiv preprint arXiv:1910.12249*, 2019.
- 329 [14] Evans, D., Nguyen-Tuong, A., and Knight, J. Effectiveness of moving target defenses. In *Moving target  
330 defense*, pp. 29–48. Springer, 2011.
- 331 [15] Gilmer, J., Ford, N., Carlini, N., and Cubuk, E. Adversarial examples are a natural consequence of test  
332 error in noise. In *ICML*, 2019.
- 333 [16] Gomes, R., Krause, A., and Perona, P. Discriminative clustering by regularized information maximization.  
334 In *NeurIPS*, 2010.
- 335 [17] Goodfellow, I. A research agenda: Dynamic models to defend against correlated attacks. *arXiv preprint  
336 arXiv:1903.06293*, 2019.
- 337 [18] Goodfellow, I. J., Shlens, J., and Szegedy, C. Explaining and harnessing adversarial examples. *arXiv  
338 preprint arXiv:1412.6572*, 2014.
- 339 [19] Graves, A. Adaptive computation time for recurrent neural networks. *arXiv preprint arXiv:1603.08983*,  
340 2016.
- 341 [20] Guo, C., Rana, M., Cisse, M., and van der Maaten, L. Countering adversarial images using input  
342 transformations. In *ICLR*, 2018.
- 343 [21] He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learning for image recognition. In *CVPR*, 2016.
- 344 [22] Hendrycks, D. and Dietterich, T. Benchmarking neural network robustness to common corruptions and  
345 perturbations. In *ICLR*, 2019.
- 346 [23] Hill, M., Mitchell, J. C., and Zhu, S.-C. Stochastic security: Adversarial defense using long-run dynamics  
347 of energy-based models. In *ICLR*, 2021. URL <https://openreview.net/forum?id=gwFTuzzjW0>.
- 348 [24] Hoffman, J., Tzeng, E., Park, T., Zhu, J.-Y., Isola, P., Saenko, K., Efros, A., and Darrell, T. Cycada:  
349 Cycle-consistent adversarial domain adaptation. In *ICML*, 2018.
- 350

- 351 [25] Ioffe, S. and Szegedy, C. Batch normalization: Accelerating deep network training by reducing internal  
352 covariate shift. In *ICML*, 2015.
- 353 [26] Kingma, D. and Ba, J. Adam: A method for stochastic optimization. In *ICLR*, 2015.
- 354 [27] Krizhevsky, A. Learning multiple layers of features from tiny images. Technical report, University of  
355 Toronto, 2009.
- 356 [28] Liang, J., Hu, D., and Feng, J. Do we really need to access the source data? source hypothesis transfer for  
357 unsupervised domain adaptation. In *ICML*, 2020.
- 358 [29] Liang, J., Hu, D., Wang, Y., He, R., and Feng, J. Source data-absent unsupervised domain adaptation  
359 through hypothesis transfer and labeling transfer. *arXiv preprint arXiv:2012.07297*, 2020.
- 360 [30] Madry, A., Makelov, A., Schmidt, L., Tsipras, D., and Vladu, A. Towards deep learning models resistant  
361 to adversarial attacks. In *ICLR*, 2018.
- 362 [31] Olston, C., Fiedel, N., Gorovoy, K., Harmsen, J., Lao, L., Li, F., Rajashekhar, V., Ramesh, S., and Soyke, J.  
363 Tensorflow-serving: Flexible, high-performance ml serving. In *NeurIPS Workshop*, 2017.
- 364 [32] Pang\*, T., Xu\*, K., and Zhu, J. Mixup inference: Better exploiting mixup to defend adversarial attacks. In  
365 *ICLR*, 2020.
- 366 [33] Perez, E., Strub, F., De Vries, H., Dumoulin, V., and Courville, A. Film: Visual reasoning with a general  
367 conditioning layer. In *AAAI*, 2018.
- 368 [34] Quionero-Candela, J., Sugiyama, M., Schwaighofer, A., and Lawrence, N. D. *Dataset shift in machine  
369 learning*. The MIT Press, 2009.
- 370 [35] Radosavovic, I., Johnson, J., Xie, S., Lo, W.-Y., and Dollár, P. On network design spaces for visual  
371 recognition. In *ICCV*, 2019.
- 372 [36] Radosavovic, I., Kosaraju, R. P., Girshick, R., He, K., and Dollár, P. Designing network design spaces. In  
373 *CVPR*, 2020.
- 374 [37] Raff, E., Sylvester, J., Forsyth, S., and McLean, M. Barrage of random transforms for adversarially robust  
375 defense. In *CVPR*, 2019.
- 376 [38] Rice, L., Wong, E., and Kolter, Z. Overfitting in adversarially robust deep learning. In *ICML*, 2020.
- 377 [39] Rony, J., Hafemann, L. G., Oliveira, L. S., Ayed, I. B., Sabourin, R., and Granger, E. Decoupling direction  
378 and norm for efficient gradient-based l2 adversarial attacks and defenses. In *CVPR*, 2019.
- 379 [40] Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A.,  
380 Bernstein, M., et al. ImageNet large scale visual recognition challenge. *IJCV*, 2015.
- 381 [41] Saenko, K., Kulis, B., Fritz, M., and Darrell, T. Adapting visual category models to new domains. In  
382 *ECCV*, 2010.
- 383 [42] Samangouei, P., Kabkab, M., and Chellappa, R. Defense-GAN: Protecting classifiers against adversarial  
384 attacks using generative models. In *ICLR*, 2018.
- 385 [43] Schneider, S., Rusak, E., Eck, L., Bringmann, O., Brendel, W., and Bethge, M. Improving robustness  
386 against common corruptions by covariate shift adaptation. In *NeurIPS*, 2020.
- 387 [44] Sehwal, V., Mahloujifar, S., Handina, T., Dai, S., Xiang, C., Chiang, M., and Mittal, P. Improving  
388 adversarial robustness using proxy distributions. *arXiv preprint arXiv:2104.09425*, 2021.
- 389 [45] Shannon, C. A mathematical theory of communication. *Bell system technical journal*, 27, 1948.
- 390 [46] Sharma, Y. and Chen, P.-Y. Attacking the madry defense model with  $l_1$ -based adversarial examples.  
391 *arXiv preprint arXiv:1710.10733*, 2017.
- 392 [47] Shazeer, N., Mirhoseini, A., Maziarz, K., Davis, A., Le, Q., Hinton, G., and Dean, J. Outrageously large  
393 neural networks: The sparsely-gated mixture-of-experts layer. In *ICLR*, 2017.
- 394 [48] Shelhamer, E., Wang, D., and Darrell, T. Blurring the line between structure and learning to optimize and  
395 adapt receptive fields. *arXiv preprint arXiv:1904.11487*, 2019.
- 396 [49] Shi, C., Holtz, C., and Mishne, G. Online adversarial purification based on self-supervised learning. In  
397 *ICLR*, 2021. URL [https://openreview.net/forum?id=\\_i3ASpP12WS](https://openreview.net/forum?id=_i3ASpP12WS).

- 398 [50] Song, Y., Kim, T., Nowozin, S., Ermon, S., and Kushman, N. Pixeldefend: Leveraging generative models  
399 to understand and defend against adversarial examples. In *ICLR*, 2018.
- 400 [51] Stutz, D., Hein, M., and Schiele, B. Disentangling adversarial robustness and generalization. In *CVPR*,  
401 June 2019.
- 402 [52] Su, D., Zhang, H., Chen, H., Yi, J., Chen, P.-Y., and Gao, Y. Is robustness the cost of accuracy?—a  
403 comprehensive study on the robustness of 18 deep image classification models. In *ECCV*, 2018.
- 404 [53] Sun, Y., Wang, X., Liu, Z., Miller, J., Efros, A. A., and Hardt, M. Test-time training for out-of-distribution  
405 generalization. In *ICML*, 2020.
- 406 [54] Szegedy, C., Zaremba, W., Sutskever, I., Bruna, J., Erhan, D., Goodfellow, I., and Fergus, R. Intriguing  
407 properties of neural networks. In *ICLR*, 2014.
- 408 [55] Tramer, F. and Boneh, D. Adversarial training and robustness for multiple perturbations. In *NeurIPS*, 2019.
- 409 [56] Tramer, F., Carlini, N., Brendel, W., and Madry, A. On adaptive attacks to adversarial example defenses.  
410 In *NeurIPS*, 2020.
- 411 [57] Veit, A. and Belongie, S. Convolutional networks with adaptive inference graphs. In *ECCV*, 2018.
- 412 [58] Wang, D., Shelhamer, E., Liu, S., Olshausen, B., and Darrell, T. Fully test-time adaptation by entropy  
413 minimization. In *ICLR*, 2021.
- 414 [59] Wang, X., Yu, F., Dou, Z.-Y., Darrell, T., and Gonzalez, J. E. Skipnet: Learning dynamic routing in  
415 convolutional networks. In *ECCV*, 2018.
- 416 [60] Wong, E., Rice, L., and Kolter, J. Z. Fast is better than free: Revisiting adversarial training. In *ICLR*, 2020.
- 417 [61] Wu, D., Xia, S.-T., and Wang, Y. Adversarial weight perturbation helps robust generalization. In *NeurIPS*,  
418 2020.
- 419 [62] Yang, B., Bender, G., Le, Q. V., and Ngiam, J. Condconv: Conditionally parameterized convolutions for  
420 efficient inference. In *NeurIPS*, 2019.
- 421 [63] Yuan, X., He, P., Zhu, Q., and Li, X. Adversarial examples: Attacks and defenses for deep learning.  
422 *TNNLS*, 2019.
- 423 [64] Zagoruyko, S. and Komodakis, N. Wide residual networks. *arXiv preprint arXiv:1605.07146*, 2016.
- 424 [65] Zhang, H., Yu, Y., Jiao, J., Xing, E., El Ghaoui, L., and Jordan, M. I. Theoretically principled trade-off  
425 between robustness and accuracy. In *ICML*, 2019.
- 426 [66] Zhang, H., Chen, H., Xiao, C., Goyal, S., Stanforth, R., Li, B., Boning, D., and Hsieh, C.-J. Towards  
427 stable and efficient training of verifiably robust neural networks. In *ICLR*, 2020.

## 428 Checklist

- 429 1. For all authors...
- 430 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s  
431 contributions and scope? **[Yes]** 1. The last move advantage is described in Section 3.2  
432 under “Discussion”. 2. Dent is unique in its adaptation of the model parameters during  
433 testing without altering training. Section 2 explains that generative and self-supervised  
434 defenses optimize the input, but require auxiliary models, and their parameters are fixed.  
435 3. Table 1 is our main result showing that dent improves state-of-the-art adversarial  
436 training defenses. 4. Adaptive attacks are explained and reported in Section 4.5.
- 437 (b) Did you describe the limitations of your work? **[Yes]** Section 4.2 describes the threat  
438 model for our work and its requirement of test batches. This is discussed more in  
439 Section 5.
- 440 (c) Did you discuss any potential negative societal impacts of your work? **[No]** Our project  
441 is about defense, not attack, and so it has less potential for direct negative impact.
- 442 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
443 them? **[Yes]**

- 444 2. If you are including theoretical results...
- 445 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 446 (b) Did you include complete proofs of all theoretical results? [N/A]
- 447 3. If you ran experiments...
- 448 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
- 449 mental results (either in the supplemental material or as a URL)? [Yes] The experiment
- 450 setup (Section 4.1) and code in the supplementary material specify experiment details
- 451 and reproduce our main results for dent+.
- 452 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were
- 453 chosen)? [Yes] Our method alters testing, not training, but its test-time optimization
- 454 details are given by the above.
- 455 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
- 456 ments multiple times)? [No] We evaluate with the standard adversarial benchmarks or
- 457 AutoAttack and RobustBench, which do not report variability in this way.
- 458 (d) Did you include the total amount of compute and the type of resources used (e.g., type
- 459 of GPUs, internal cluster, or cloud provider)? [No] Our project does not require exotic
- 460 resources, but standard deep learning hardware. In particular, we use V100 and Titan
- 461 Xp GPUs on local servers.
- 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] Libraries, datasets,
- 464 and models are all cited in Experiments (Section 4.
- 465 (b) Did you mention the license of the assets? [No] The licenses are mentioned at the given
- 466 references and URLs.
- 467 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
- 468 Our code is included in the supplemental material.
- 469 (d) Did you discuss whether and how consent was obtained from people whose data you're
- 470 using/curating? [N/A] We did not collect new data.
- 471 (e) Did you discuss whether the data you are using/curating contains personally identifiable
- 472 information or offensive content? [N/A] We did not collect new data.
- 473 5. If you used crowdsourcing or conducted research with human subjects...
- 474 (a) Did you include the full text of instructions given to participants and screenshots, if
- 475 applicable? [N/A]
- 476 (b) Did you describe any potential participant risks, with links to Institutional Review
- 477 Board (IRB) approvals, if applicable? [N/A]
- 478 (c) Did you include the estimated hourly wage paid to participants and the total amount
- 479 spent on participant compensation? [N/A]