

Why certified defences hurt robust generalisation

Anonymous Authors¹

Abstract

Several works have shown that state-of-the-art classifiers are vulnerable to adversarial examples, raising serious concerns for their deployment in safety-critical applications. To address this issue two classes of methods have emerged: empirical defences and certified defences. While the former has good robustness, but no guarantees, the latter sacrifices some robustness in exchange for guarantees. Until now, these two approaches have not been systematically compared in the literature, and a clear understanding of the robustness gap remains elusive. In this paper, we show through extensive empirical evidence that models trained with certified defences suffer from worse accuracy, robustness and fairness than empirical defences. Further, we identify three key factors contributing to the robustness gap between the two approaches and support our arguments with both theoretical and experimental evidence. We hope this serves as a guide to practitioners regarding the use of various defences as well as a motivation for designing future methods.

1. Introduction

Several works have shown that state-of-the-art classifiers are vulnerable to adversarial examples, i.e., imperceptible perturbations to the input that can change the classifier prediction (Biggio et al., 2013; Szegedy et al., 2014). Hence, robustness to adversarial examples has become a crucial design goal when deploying machine learning models in safety critical applications. In real-world scenarios, robustness against many different types of input perturbations may be desired depending on the domain of application. In this paper, we consider the well-studied ℓ_p -ball threat model, where $\mathcal{B}_\epsilon := \{\delta : \|\delta\|_p \leq \epsilon\}$ represents the set of allowed perturbations for some ℓ_p -ball with radius ϵ centred

around the origin. Then, for any distribution \mathcal{D} , classifier $f_\theta : \mathbb{R}^d \rightarrow \mathbb{R}^k$ parameterised by $\theta \in \mathbb{R}^\rho$, and loss function L , adversarial defences aim to solve the following robust optimisation problem:

$$\min_{\theta} \mathbf{R}_\epsilon(\theta) := \mathbb{E}_{(x,y) \sim \mathcal{D}} \left[\max_{\delta \in \mathcal{B}_\epsilon} L(f_\theta(x + \delta), y) \right] \quad (1)$$

We call $\mathbf{R}_\epsilon(\theta)$ the robust error when L is the 0-1 loss function. In practice, as the distribution \mathcal{D} is unknown, we minimise the empirical robust error on a finite dataset D sampled from \mathcal{D} . Further, in the case of neural networks, the inner-maximisation is a non-convex optimisation problem and prohibitively hard to solve from a computational perspective (Katz et al., 2017; Weng et al., 2018). Instead, two efficient techniques are widely used to overcome the computational barrier: *empirical* defences that provide a lower bound and *certified* defences that provide an upper bound on the solution.

Adversarial training (AT) (Goodfellow et al., 2015; Madry et al., 2018) is one of the most popular empirical defences to date. AT minimises the worst-case empirical loss in Equation (1) by approximately solving the inner-maximisation problem with first-order optimisation methods. However, despite its simplicity and computational efficiency, AT does not provide any robustness guarantees, which are essential in many safety-critical domains.

To address this limitation, there has been significant interest in designing certified defences, i.e., methods for learning neural networks that are *provably* robust to ℓ_p -ball perturbations on the test data, allowing a safe upper bound estimate on the robust error. Many recent works have proposed to solve a convex relaxation of the inner-maximisation problem by relaxing the non-convex ReLU constraint sets with convex ones (Wong & Kolter, 2018; Raghunathan et al., 2018; DVjiotham et al., 2018; Zhang et al., 2020). Nevertheless, certified defences based on convex relaxations suffer from an inherent flaw: the upper bound they provide on the robust error is far from being tight (Salman et al., 2019).

While empirical defences have good robustness, but no guarantees, certified defences appear to sacrifice some robustness in exchange for guarantees. However, to the best of our knowledge, a systematic comparison and understanding of the gap in robustness between the two classes of methods is largely missing in the literature.

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

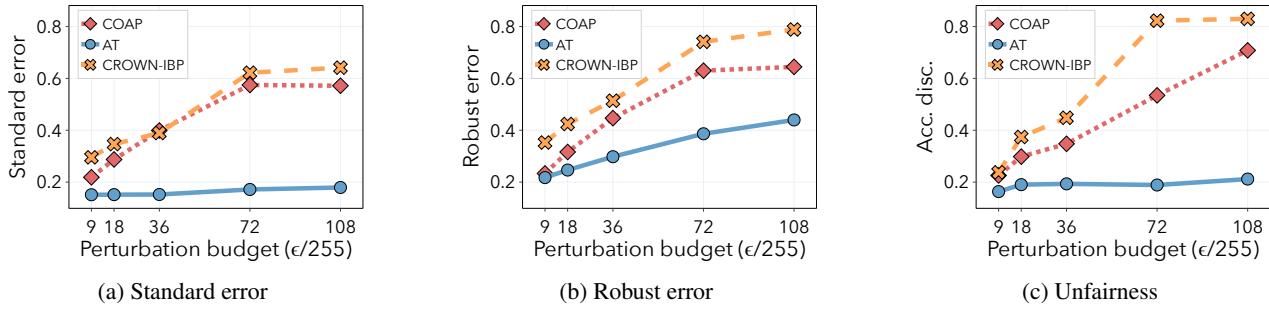


Figure 1. Results for ℓ_2 -ball perturbations on the CIFAR-10 test set. We compare ResNet architectures trained using state-of-the-art certified defences CROWN-IBP (Zhang et al., 2020; Xu et al., 2020) and COAP (Wong et al., 2018; Wong & Kolter, 2018) against the most popular empirical defence to date AT (Goodfellow et al., 2015; Madry et al., 2018). In Figures 1a and 1b, we plot the standard error and robust error respectively as the perturbation budget increases. In Figure 1c, we plot accuracy discrepancy, i.e. the difference between overall accuracy and worst class accuracy, as the perturbation budget increases. We refer the reader to Section 2 for further details on the models and robust evaluation.

In this paper, we provide the first systematic comparison between empirical and certified defences on popular computer vision datasets. Further, we identify three key factors contributing to the robustness gap between the two approaches and support our arguments with both theoretical and experimental evidence. We hope this serves as a guide to practitioners regarding the use of various defences as well as a motivation for designing future defence methods.

Specifically, our contributions are as follows:

- In Section 2, we show through extensive experimental evidence on image datasets that models trained with certified defences suffer from worse accuracy and robustness than empirical defences (see Figures 1a and 1b).
- In Section 3, we identify three factors driving the robustness gap: (i) the number of active neurons; (ii) the magnitude of the adversarial perturbations compared to the implicit margin of the data; (iii) the alignment between the adversarial perturbations and the “signal” direction. We provide both theoretical and experimental evidence on synthetic and image datasets to motivate these findings.
- In Section 4, we investigate the fairness of empirical and certified defences; we show that certified defences lead to models with worse fairness than empirical defences (see Figure 1c).

1.1. Related work

Deterministic defences A large body of literature has been devoted to developing methods for training neural networks that are provably robust to ℓ_p -ball perturbations. These works include methods based on semidefinite relaxations (Raghunathan et al., 2018), convex relaxations

(Wong & Kolter, 2018; Wong et al., 2018), abstract interpretation (Mirman et al., 2018; Singh et al., 2018), and interval bound propagation (Gowal et al., 2019a; Zhang et al., 2020). Salman et al. (2019) unify these different views in a common convex relaxation framework and show that convex relaxations that approximate each ReLU output separately suffer an inherent tightness barrier. In particular, even the optimal convex relaxation cannot obtain tight bounds on the robust error. However, Salman et al. (2019) investigate the tightness of convex relaxations for verification purposes, i.e., for certifying the robustness of already trained models. In contrast, we explore the impact of this looseness on the robustness of models trained using convex relaxations, which to the best of our knowledge has not been systematically studied before. Previous research in this direction, such as the work of Jovanovic et al. (2021), compares the training dynamics of various convex relaxations. Our work, on the other hand, focuses on the comparison between empirical and certified defences.

Probabilistic defences As an alternative to deterministic defences, probabilistic defences give a guarantee of robustness with a certain probability. One of the most popular probabilistic defence is randomised smoothing (Cohen et al., 2019; Lécuyer et al., 2019; Li et al., 2019). The key idea behind this technique is to transform an arbitrary classifier $f_\theta : \mathbb{R}^d \rightarrow \mathbb{R}^k$ into a “smoothed” classifier $g_\theta : \mathbb{R}^d \rightarrow \mathbb{R}^k$. In particular, for a given data point $x \in \mathbb{R}^d$ and variance σ^2 , the smoothed classifier’s prediction $g_\theta(x)$ is defined as the most probable prediction by f_θ on the random variable $z \sim \mathcal{N}(x, \sigma^2 I_d)$. Despite its popularity, randomised smoothing significantly increases the computational cost of inference and suffers from several limitations. For example, Mohapatra et al. (2021) investigate the side-effects of randomised smoothing and show that it significantly hurts the disparity in class-wise accuracy; Sun et al. (2022) ob-

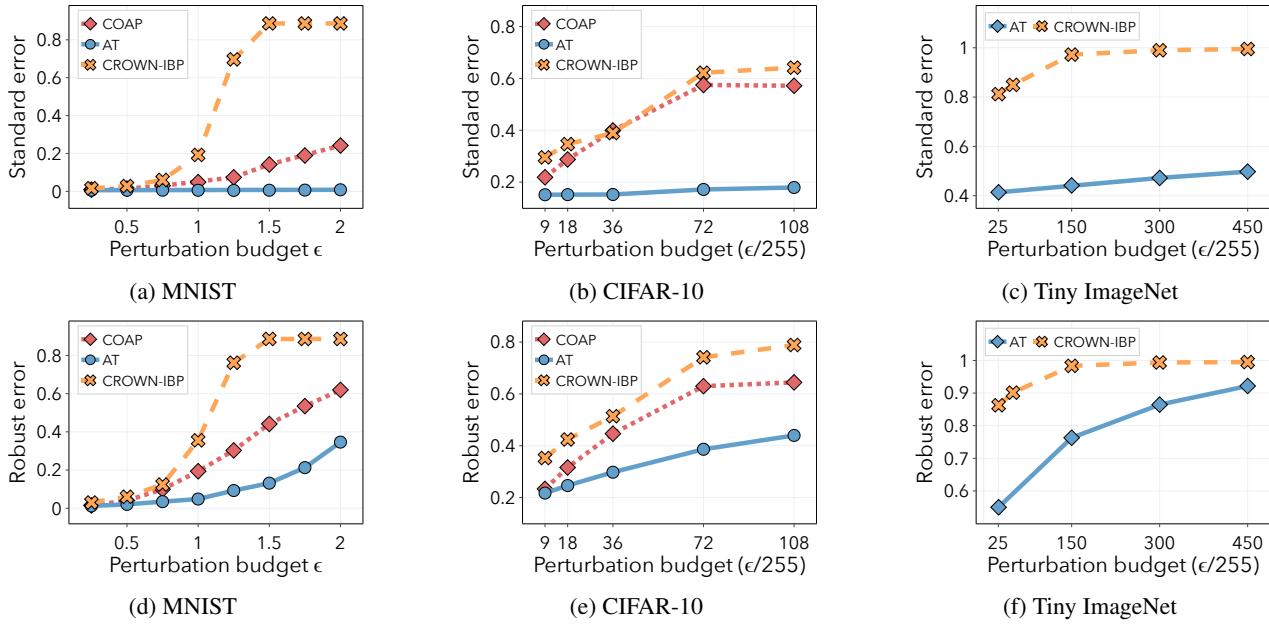


Figure 2. Results for ℓ_2 -ball perturbations on MNIST, CIFAR-10 and Tiny ImageNet test sets. We compare ResNet architectures trained using state-of-the-art certified defences CROWN-IBP (Zhang et al., 2020; Xu et al., 2020) and COAP (Wong et al., 2018; Wong & Kolter, 2018) against the most popular empirical defence to date AT (Goodfellow et al., 2015; Madry et al., 2018). In Figures 2a to 2c we plot the standard error as the perturbation budget increases. In Figures 2d to 2f we plot the robust error as the perturbation budget increases. We omit COAP in Figures 2c and 2f as it does not scale to Tiny ImageNet.

served that randomised smoothing is extremely vulnerable to low-frequency corruptions of the test data. Further, several works have exposed an accuracy-robustness tradeoff related to the smoothness of the classifier (Yang et al., 2020; Blum et al., 2020; Kumar et al., 2020). Compared to convex relaxation-based methods, the drawbacks of randomised smoothing are much better understood, and efforts are being made towards developing new defences to bridge the gap with empirical defences (Nandy et al., 2022). Therefore, we focus our attention in this paper on understanding the limitations of convex relaxation-based methods.

2. Generalisation gap between empirical and certified defences

In this section, we provide a systematic comparison of empirical and certified defences on three real-world computer vision datasets: Tiny ImageNet (Le & Yang, 2015), CIFAR-10 (Krizhevsky, 2009) and MNIST (LeCun et al., 1998). This paper mainly focuses on ℓ_2 -ball perturbations. Supplementary experiments with ℓ_∞ -ball perturbations, presented in Appendix E, indicate a similar trend. Moreover, we note that certified defences do not scale to larger vision datasets such as ImageNet, hence why we omit them in our experiments.

Methods Among certified defences based on convex relaxations, we consider the convex outer adversarial polytope (COAP) (Wong et al., 2018; Wong & Kolter, 2018), which achieves state-of-the-art certified robustness under ℓ_2 -ball perturbations. We also consider CROWN-IBP (Zhang et al., 2018; Xu et al., 2020), which combines the tight convex relaxation CROWN (Zhang et al., 2018) with interval bound propagation (IBP) (Gowal et al., 2019b; Mirman et al., 2018) and achieves state-of-the-art certified robustness under ℓ_∞ -ball perturbations. We compare these methods against the most popular empirical defence—adversarial training (AT) (Goodfellow et al., 2015; Madry et al., 2018), which is the de-factor method for adversarial robustness among practitioners.

Models and robust evaluation To reliably evaluate the defences in the ℓ_2 -ball perturbations threat model, we use the strongest version of AutoAttack (AA+) (Croce & Hein, 2020). For CIFAR-10, we train a residual network (ResNet) and for MNIST we train a convolutional neural network (CNN). Both architectures were introduced in Wong et al. (2018) as standard benchmarks for certified defences. For Tiny ImageNet, we train a WideResNet along the same lines of Xu et al. (2020). We refer the reader to Appendix D.4 for complete experimental details.

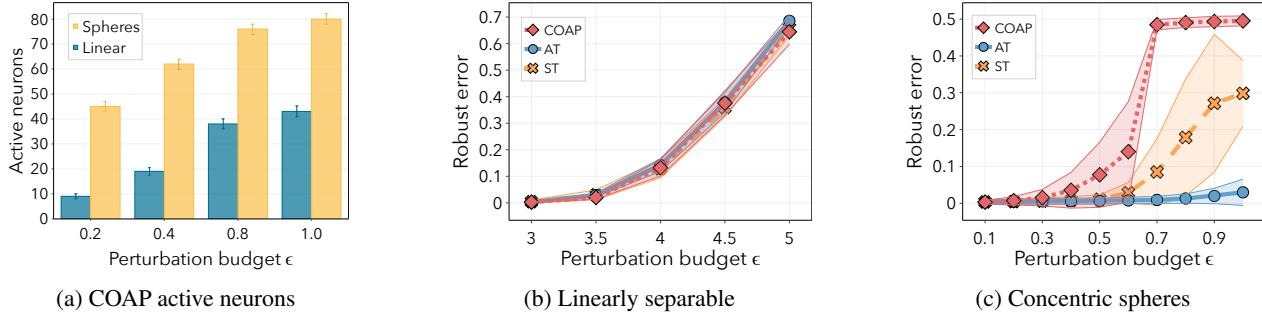


Figure 3. We report mean and standard deviation over 15 seeds. In Figure 3a, we plot the number of neurons in the activation set for the linearly separable and concentric spheres distribution after they have been trained. In Figures 3b and 3c we plot the robust error for standard training (ST), adversarial training (AT) and convex outer adversarial polytope (COAP), when training on the linearly separable and concentric spheres distributions respectively. See Appendix D.2 for complete experimental details.

Results Several studies have shown that adversarial training may lead to an increase in standard error when compared with standard training (Raghunathan et al., 2020; Tsipras et al., 2019; Zhang et al., 2019). We observe the same phenomenon to a much higher degree when using certified defences. In particular, our experimental results show that certified defences not only suffer worse standard error but also worse robust error than adversarial training. Figures 2a and 2b shows that for increasing perturbation budget, the standard error gap between certified (CROWN-IBP, COAP) and empirical defences (AT) increases for both MNIST and CIFAR-10 respectively. Specifically, the gap reaches almost 90% for CROWN-IBP on MNIST when $\epsilon = 1.5$. Figures 2d and 2e show that the robust error gap also increases with increasing perturbation budget for both MNIST and CIFAR-10. In particular, the gap reaches almost 40% for the largest perturbation budgets. Additionally, we observe a significant standard and robust error gap between AT and CROWN-IBP for Tiny ImageNet in Figures 2c and 2f.

Our experiments with popular image datasets and neural network architectures reveal a significant standard and robust generalisation gap between certified and empirical defences. We will explain the reasons for this behaviour, both theoretically and experimentally, in the following section.

3. Underlying factors of the generalisation gap

In this section, we investigate when certified defences perform poorly in practice. For the sake of clarity, we focus our attention on COAP, as it showed better robustness than CROWN-IBP when tested on image datasets in Section 2. Nevertheless, we expect that the takeaways from this section will transfer to CROWN-IBP as well.

In particular, we identify three factors underlying the robust generalisation gap: (i) the number of active neurons; (ii) the magnitude of the adversarial perturbations compared to

the implicit margin of the data; (iii) the alignment between the adversarial perturbations and the signal direction. We verify these factors using image datasets and two synthetic data distributions: a linearly separable distribution similar to those studied Clarysse et al. (2022); Tsipras et al. (2019), and the concentric spheres distribution studied in Gilmer et al. (2018); Nagarajan & Kolter (2019).

Data and threat models Similar to the previous section, we focus on ℓ_2 -ball perturbations of radius ϵ . To sample a data point from the linearly separable distribution with margin $\gamma > 0$, first we sample the label $y \in \{+1, -1\}$ with equal probability. Then, sample $\tilde{x} \in \mathbb{R}^{d-1}$ from a standard normal distribution $\tilde{x} \sim \mathcal{N}(0, \sigma^2 I_{d-1})$ and set the covariate vector $x = [\gamma \text{sgn}(y); \tilde{x}]$, where $[;]$ is the concatenation operator. To sample from the concentric spheres distribution with radii $0 < R_1 < R_{-1}$, first draw a binary label $y \in \{+1, -1\}$ with equal probability and then the covariate vector $x \in \mathbb{R}^d$ is distributed uniformly on the sphere of radius R_y . Observe that achieving a low test error on the concentric spheres distribution requires a non-linear classifier.

3.1. Factor (i): Active neurons

The first factor we investigate is the number of active neurons for models trained with certified defences. The intuition behind this factor is the following: COAP relaxes the non-convex ReLU constraints for all neurons that activate within the perturbation set, i.e., when there exists a perturbation $\delta \in \mathcal{B}_\epsilon$ for which the input to the neuron crosses 0. Hence, the larger the number of active neurons, the worse the approximation. In Figure 3a we plot the number of active neurons on the concentric spheres and linear distributions as the perturbation budget increases. Observe that the amount is much higher for the concentric spheres than for the linearly separable distribution. This is consistent with the intuition that the complex spherical decision bound-

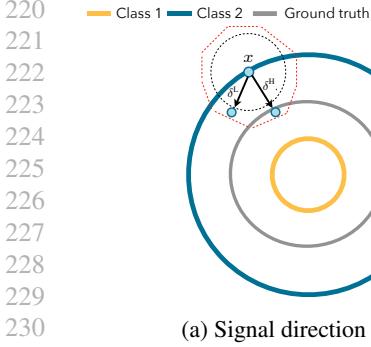
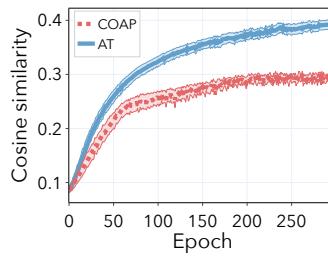
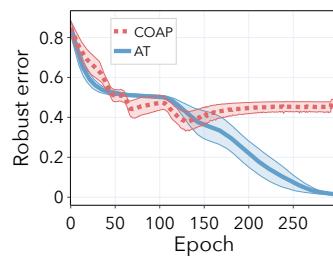

 220
221
222
223
224
225
226
227
228
229
230
231
(a) Signal direction

 232
233
234
235
236
237
238
239
240
241
(b) Signal alignment

 232
233
234
235
236
237
238
239
240
241
(c) Robust error

Figure 4. We report mean and standard error over 15 seeds. In Figure 4a we plot a 2-dimensional visualisation of the concentric spheres dataset; the black dashed circle represents the ℓ_2 -ball threat model and the red dashed polytope illustrates the convex approximation. Observe that for perturbations that are not significantly aligned with the signal direction, e.g. δ^L , the perturbed data point cannot cross the decision boundary. On the other hand, for perturbations that are aligned with the signal direction, e.g. δ^H , which overlaps with the signal direction, the perturbed data points can cross the decision boundary. In Figure 4b we plot the average cosine similarity between ℓ_2 norm-bounded perturbations on the training data and the signal directed vector. In Figure 4c we plot the robust error for adversarial training (AT) and convex outer adversarial polytope (COAP), for large perturbation budget $\epsilon = 1.0$. We observe that when cosine similarity is high, the gap in robust error between COAP and AT increases. Hence, especially approximations in the signal direction can hurt robust generalisation.

ary requires much more active neurons compared to the linear decision boundary which only needs 1. Moreover, in Figures 3b and 3c, we plot the robust error of standard training (ST), adversarial training (AT), and convex outer adversarial polytope (COAP) on the two distributions. We see that in contrast to the linear setting, COAP has a much higher robust error on the concentric spheres distribution than AT and ST. Overall, these findings indicate that the number of active neurons is a key factor underlying the robust generalisation gap between certified and empirical defences.

3.2. Factor (ii): Alignment between the adversarial perturbations and the signal direction

The second factor we investigate is the alignment between adversarial perturbations and the “signal” direction. We define the signal direction $s(x, y)$ for a data point (x, y) and a ground truth f^* as the direction along which we can flip the label with minimal perturbation budget. More formally,

$$s(x, y) = \underset{\delta}{\operatorname{argmin}} \min_{\text{s.t. } f^*(x + \epsilon \cdot \delta) \cdot y < 0} \epsilon \quad (2)$$

For a point (x, y) drawn from the concentric spheres distribution, the signal direction is given by $s(x) = y \frac{x}{\|x\|_2}$. Observe that, the convex approximation essentially magnifies the radius of the ℓ_2 -ball perturbation, and thus, a significant alignment of the ℓ_2 -ball perturbation with the signal direction may result in points crossing the ground truth decision boundary, as illustrated in Figure 4a.

We support this argument with experiments on the concentric spheres dataset. In Figure 4b we plot the cosine

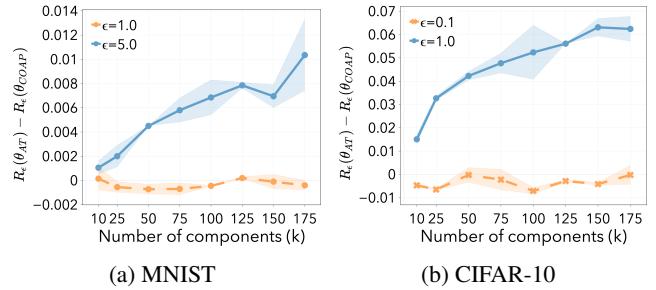


Figure 5. We report mean and standard deviation over 2 seeds. In Figures 5a and 5b we plot the robust error gap between AT and COAP as the threat model captures more signal direction. We omit Tiny ImageNet as COAP cannot scale to datasets of large size. We refer the reader to Appendix D.5 for complete experimental details.

similarity between the ℓ_2 -ball perturbations computed on the training set and the signal direction, for large perturbation budget. By comparing the signal alignment in figure 4b with the corresponding robust test error in Figure 4c, we can see that during the early stages of training, the ℓ_2 -ball perturbations are not aligned with the signal direction and the robust error for COAP is similar to AT. However, as training progresses and the perturbations begin to align with the signal direction, the robust error gap between COAP and AT increases. This provides evidence that as training progresses, ℓ_2 -ball perturbations become significantly aligned with the signal direction and the robustness gap worsens.

So far, we have intentionally investigated the factors on very controlled settings to gain a better understanding. Now, we

275 experiment with more realistic datasets, namely MNIST
 276 and CIFAR-10. As it is difficult to model the exact signal
 277 direction for image datasets, we propose an approximation
 278 of the signal-directed threat model introduced in Section 3.3.
 279 Using this, we demonstrate how the robust generalisation
 280 gap between AT and COAP gradually increases as the threat
 281 model becomes more closely aligned with the signal direc-
 282 tion and the perturbation budget ϵ increases.

283 Below, we describe the approximation for the signal-
 284 directed threat model. First, randomly sample k unit vectors
 285 s_k and define the threat model as follows:

$$\mathcal{B}_\epsilon(x) = \cup_k \{z_1 = x + s_k \beta \mid |\beta| \leq \epsilon\} \quad (3)$$

287 As k increases, so does the probability that the signal
 288 direction will be significantly aligned with at least one of
 289 the k unit vectors. Then, we train both COAP and AT to
 290 be robust under this threat model. In Figures 5a and 5b
 291 we observe that the robust error gap for large perturbation
 292 budget between AT and COAP significantly increases as k
 293 increases, for both MNIST and CIFAR-10. This provides
 294 further evidence that *perturbations aligned with the signal*
 295 *direction worsen the robustness gap.*

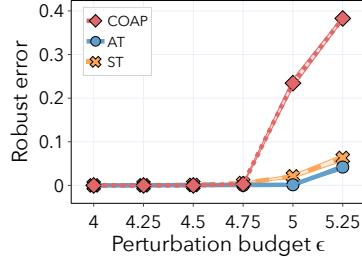
296 Further, the robust generalisation gap of MNIST is signif-
 297 icantly lower than that of CIFAR-10. This is likely due to
 298 the fact that the MNIST dataset is almost linearly separable,
 299 meaning that the decision boundary between classes
 300 is relatively simple and can be approximated by an hyper-
 301 plane. We already displayed a similar distinction between
 302 the linearly separable and the concentric spheres distribu-
 303 tions in Figure 3a. In light of these observations, it is not
 304 surprising that the robustness gap is less pronounced for
 305 MNIST than for CIFAR-10.

306 Finally, note how the robust generalisation gap is only ob-
 307 served for adversarial perturbations with large magnitude.
 308 This quantity plays an important role for robust generalisation,
 309 as shown already in Figure 3c, where the gap in robust
 310 error between certified and empirical defences increases
 311 with larger magnitudes.

312 3.3. Factor (iii): Magnitude of the adversarial 313 perturbations

314 The third factor we investigate is the magnitude of the adver-
 315 sarial perturbations. In this section, we present a theoretical
 316 result which determines the minimum magnitude required
 317 to observe the robust generalisation gap. Further, we corrob-
 318 orate our theoretical result with experimental evidence on
 319 synthetic data.

320 **Data and threat model** We consider the linearly separable
 321 distribution. Note that for this distribution we did not
 322 observe a robust generalisation gap in Figure 3b. Since the



283 *Figure 6.* We report mean and standard deviation over 15 seeds.
 284 We consider the linearly separable distribution and the signal-
 285 directed threat model defined in Equation (4). We plot the robust
 286 error for standard training (ST), adversarial training (AT) and
 287 convex outer adversarial polytope (COAP) as the perturbation
 288 budget ϵ increases. See Appendix D.1 for complete experimental
 289 details.

329 alignment with the signal direction can increase the robust-
 330 ness gap – as suggested in Section 3.2 – we choose our threat
 331 model to maximise the alignment between perturbations and
 332 signal direction. In fact, we consider perturbations that are
 333 perfectly aligned with the signal direction. Observe that
 334 the signal direction, as defined in Equation (2), corresponds
 335 to the first index of the input. Hence, we define the set of
 336 allowed perturbations

$$\mathcal{B}_\epsilon(x) = \{z_1 = x + e_1 \beta \mid |\beta| \leq \epsilon\}, \quad (4)$$

337 where e_1 is the canonical basis vector of the first coordi-
 338 nate. We refer the reader to Appendix A.1, where we
 339 derive an extension of the convex outer adversarial poly-
 340 tope (COAP) (Wong & Kolter, 2018) to this setting.

341 We are now ready to present our theoretical result for the
 342 linearly separable distribution. Perhaps surprisingly, The-
 343 orem 1 shows for a simple neural network that, in high
 344 dimensions, certified defences (COAP) yield higher robust
 345 error than empirical defences (AT) when the perturbation
 346 budget is large. Below, we outline our theoretical setting in
 347 more detail.

348 **One-step gradient descent** We study the early phase of
 349 neural network optimisation. Under structural assumptions
 350 on the data, it has been proven that one gradient step with
 351 sufficiently large learning rate can drastically decrease the
 352 training loss (Chatterji et al., 2021) and extract task-relevant
 353 features (Frei et al., 2022; Daniely & Malach, 2020). A
 354 similar setting was also studied recently in Ba et al. (2022)
 355 for the MSE loss in the high-dimensional asymptotic limit.
 356 We focus on the classification setting with binary cross-
 357 entropy loss.

330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384

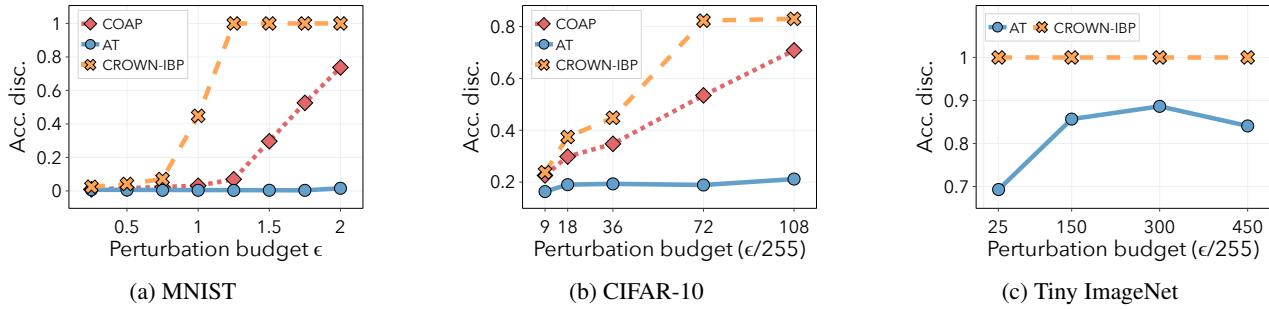


Figure 7. Results for ℓ_2 -ball perturbations on MNIST, CIFAR-10 and Tiny ImageNet test sets. In Figures 7a to 7c we plot the accuracy discrepancy as the perturbation budget increases. We refer the reader to Appendix D.4 for complete experimental details.

One-neuron neural network Consider the hypothesis class of one-neuron shallow neural networks $f_\theta : \mathbb{R}^d \rightarrow \mathbb{R}$,

$$f_\theta(x) = a \operatorname{ReLU}(\theta^\top x) + b \quad (5)$$

where $x, \theta \in \mathbb{R}^d$ and $a, b \in \mathbb{R}$ and the only trainable parameter is θ_1 . Note that as our distribution is linearly separable, our hypothesis class includes the ground truth.

Below we state our main theorem.

Theorem 1. Let \mathbf{R}_ϵ be the robust risk of f_θ , defined as

$$\mathbf{R}_\epsilon(\theta) := \mathbb{P}_{(x,y)} [\exists z \in \mathcal{B}_\epsilon(x) : y \neq \operatorname{sgn}(f_\theta(z))].$$

Let $\bar{\theta}$ and $\tilde{\theta}$ be the network parameters after one step of gradient descent with respect to AT and COAP objectives. If the initialised network parameters θ satisfy

$$\frac{\|\theta_{2:d}\|_2}{\|\theta_1\|_2} > \sqrt{\max \left(\frac{(7\epsilon - \gamma)(\gamma + \epsilon)^4}{4\sigma^2(\gamma^2 - 10\gamma\epsilon + 13\epsilon^2)}, \frac{(\gamma + \epsilon)^3}{12\sigma^2\epsilon} \right)}$$

and $\frac{2}{3}\gamma < \epsilon < \gamma$, COAP yields higher robust risk than AT:

$$\mathbf{R}_\epsilon(\tilde{\theta}) > \mathbf{R}_\epsilon(\bar{\theta}).$$

Theorem 1 relies on two main assumptions. The first is an assumption on the data dimensionality and the initialisation of the network parameters θ . For instance, if we initialise the network parameters θ by sampling from a sub-gaussian distribution, then the euclidean norm $\|\theta\|_2$ concentrates around \sqrt{d} with high probability. Hence, the assumption is satisfied when the data dimensionality d is sufficiently high. The second assumption requires that the perturbation budget ϵ is sufficiently close to the implicit margin γ of the data. This provides further evidence that a large magnitude is required to observe the generalisation gap between empirical and certified defences. Moreover, it is consistent with our experimental evidence on image datasets, as the robust generalisation gap significantly worsen for large perturbation budgets.

Synthetic experiments We corroborate our theory with experimental evidence using a one-hidden layer neural network with 100 neurons. In particular, we investigate the effect of perturbation budget ϵ on robust generalisation for three different models: standard training (ST), adversarial training (AT) and convex outer adversarial polytope (COAP). In Figure 6, we plot the robust error as the perturbation budget ϵ increases. We observe that the gap between AT and COAP increases with increasing perturbation budget, which is in agreement with Theorem 1. Overall, these findings suggest that the magnitude of the adversarial perturbations is another key factor underlying the robust generalisation gap. In the next section, we extend our investigation of the gap between certified and empirical defences to another desirable metric, namely fairness.

4. Fairness gap between certified and empirical defences

Xu et al. (2021) observed that empirical defences, such as adversarial training, may induce a large discrepancy of robustness and accuracy among different classes. However, it is not yet understood if the same holds true for certified defences. In this section we examine the fairness of both certified and empirical defences, and to the best of our knowledge, ours is the first study of the fairness gap between the two approaches. In particular, we show that certified defences lead to models with worse fairness than empirical defences. Then, we present a theoretical result which determines the minimum magnitude required to observe the fairness gap, indicating that factor (iii) affects both the robustness and fairness gaps.

We measure fairness by examining the difference between overall accuracy and worst class accuracy and refer to this metric as *accuracy discrepancy* (Buolamwini & Gebru, 2018; Sanyal et al., 2022; Xu et al., 2021). More formally, let $\mathbf{R}(\theta)$ be the standard error of the classifier f_θ and $\mathbf{R}^k(\theta)$ the standard error conditioned on the class label k . The

accuracy discrepancy is measured as follows:

$$\frac{\max_k \mathbf{R}^k(\theta) - \mathbf{R}(\theta)}{1 - \mathbf{R}(\theta)} \quad (6)$$

For MNIST, CIFAR-10 and Tiny ImageNet we observe in Figures 7a to 7c that COAP and CROWN-IBP have a significant larger accuracy discrepancy than AT.

We now present a theoretical result on fairness for the linearly separable distribution. Our theoretical setting is the same as described in Theorem 1, i.e., we study a one neuron neural network after one step of gradient descent. In Theorem 2, we show for a simple neural network that, in high dimensions, certified defences (COAP) yield higher accuracy discrepancy than empirical defences (AT) when the perturbation budget is large.

Theorem 2. *Let $\bar{\theta}$ and $\tilde{\theta}$ be the network parameters after one step of gradient descent with respect to AT and COAP objectives. Then, if the network parameters satisfy the same condition of Theorem 1 and $\frac{2}{3}\gamma < \epsilon < \gamma$, COAP yields higher accuracy discrepancy than AT:*

$$\max_k \mathbf{R}^k(\tilde{\theta}) - \mathbf{R}(\tilde{\theta}) > \max_k \mathbf{R}^k(\bar{\theta}) - \mathbf{R}(\bar{\theta}).$$

Theorem 2 relies on the same assumptions of Theorem 1. The first is an assumption on the data dimensionality which we discussed already. The second assumption requires that the perturbation budget ϵ is sufficiently close to the implicit margin γ of the data. This indicates that factor (iii) affects not only the robust generalisation gap, but also the fairness gap.

5. Discussion and future work

In this section, we discuss the implications of our observations for practitioners and for designing future defence methods. In particular, we present a simple experiment where we artificially increase the implicit margin of the data by removing some of the closest training samples to the decision boundary and observe that the robustness of COAP significantly improves.

First, we train a convolutional neural network on the entire CIFAR-10 dataset. Then, we artificially increase the distance between the classes by removing the p-percent of closest samples to the decision boundary of the trained classifier for $p \in [5, 10, 15, 20]$. Since, the distance cannot be analytically calculated for neural networks, we determine the closest samples using the confidence scores, i.e. the softmax probabilities, of the trained classifier to the true class. Hence, we remove the p-percent of points with the lowest confidence scores. In Figure 8a we plot the distribution of the confidence scores for the classifier trained on CIFAR-10, we observe that most of the scores are significantly large.

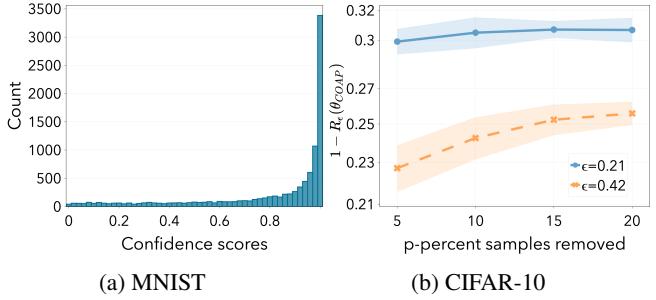


Figure 8. We report mean and standard deviation over 5 seeds. In Figure 8a we plot the distribution of the confidence scores, i.e. softmax probabilities, for a convolutional neural network classifier trained on CIFAR-10. In Figure 8b we plot COAP robust accuracy as we remove p-percent of the points with lowest confidence scores. We refer the reader to Appendix D.4 for complete experimental details.

To ensure a fair comparison for $p < 20$, we uniformly subsample the dataset such that the number of training samples is the same for all experiments.

We consider the ℓ_2 -ball perturbations threat model and we reliably evaluate the robust error using the strongest version of AutoAttack (AA+) (Croce & Hein, 2020). Figure 8b shows that the robust accuracy of COAP increases as we remove a significant percentage of samples close to the decision boundary from the training data. Further, the gain in robust accuracy is more significant when the perturbation budget is large.

The main takeaway from this experiment is that COAP suffers in robust generalisation for datasets where the implicit margin is smaller. Hence, practitioners should consider the implicit margin of the data when selecting the most appropriate defence for their specific situation. A potential approach could be using the standard accuracy as a guide, as it correlates well with the implicit margin of the data. Furthermore, future certified defences could be improved by selectively ignoring datapoints that are close to the decision boundary during the certification process. We leave both of these ideas as interesting directions for future research.

6. Conclusion

In this paper, we show that models trained with certified defences suffer from worse accuracy and robustness than empirical defences. Further, we are the first to show that certified defences also suffer from worse fairness. Finally, we develop intuition on both synthetic and image datasets for why certified defences hurt generalisation. We believe that shedding light on the robustness gap between empirical and certified defences will not only provide us with a clearer picture of the trade-offs observed in practice but also lead to better approaches for certified robustness.

References

- Ba, J., Erdogdu, M. A., Suzuki, T., Wang, Z., Wu, D., and Yang, G. High-dimensional Asymptotics of Feature Learning: How One Gradient Step Improves the Representation, 2022. arXiv:2205.01445.
- Biggio, B., Corona, I., Maiorca, D., Nelson, B., Srndic, N., Laskov, P., Giacinto, G., and Roli, F. Evasion Attacks against Machine Learning at Test Time. In *Machine Learning and Knowledge Discovery in Databases - European Conference*, 2013.
- Blum, A., Dick, T., Manoj, N., and Zhang, H. Random Smoothing Might be Unable to Certify L^∞ Robustness for High-Dimensional Images. *Journal of Machine Learning Research*, 2020.
- Buolamwini, J. and Gebru, T. Gender Shades: Intersectional Accuracy Disparities in Commercial Gender Classification. In *Proceedings of the 1st Conference on Fairness, Accountability and Transparency*, pp. 77–91. PMLR, January 2018. ISSN: 2640-3498.
- Chatterji, N. S., Long, P. M., and Bartlett, P. L. When Does Gradient Descent with Logistic Loss Find Interpolating Two-Layer Networks? *Journal of Machine Learning Research*, 2021.
- Clarysse, J., Hörrmann, J., and Yang, F. Why adversarial training can hurt robust accuracy, 2022. arXiv:2203.02006.
- Cohen, J. M., Rosenfeld, E., and Kolter, J. Z. Certified Adversarial Robustness via Randomized Smoothing. In *Proceedings of the International Conference on Machine Learning*, 2019.
- Croce, F. and Hein, M. Reliable evaluation of adversarial robustness with an ensemble of diverse parameter-free attacks. In *Proceedings of the International Conference on Machine Learning*, 2020.
- Daniely, A. and Malach, E. Learning Parities with Neural Networks. In *Advances in Neural Information Processing Systems*, 2020.
- Dvijotham, K., Stanforth, R., Gowal, S., Mann, T. A., and Kohli, P. A Dual Approach to Scalable Verification of Deep Networks. In *Proceedings of the Conference on Uncertainty in Artificial Intelligence*, 2018.
- Erdemir, E., Bickford, J., Melis, L., and Aydöre, S. Adversarial Robustness with Non-uniform Perturbations. In *Advances in Neural Information Processing Systems*, 2021.
- Frei, S., Chatterji, N. S., and Bartlett, P. L. Random Feature Amplification: Feature Learning and Generalization in Neural Networks, May 2022. arXiv:2202.07626.
- Gilmer, J., Metz, L., Faghri, F., Schoenholz, S. S., Raghu, M., Wattenberg, M., and Goodfellow, I. J. Adversarial Spheres. 2018. arXiv: 1801.02774.
- Goodfellow, I. J., Shlens, J., and Szegedy, C. Explaining and Harnessing Adversarial Examples. In *Proceedings of the International Conference on Learning Representations*, 2015.
- Gowal, S., Dvijotham, K., Stanforth, R., Bunel, R., Qin, C., Uesato, J., Arandjelovic, R., Mann, T., and Kohli, P. On the Effectiveness of Interval Bound Propagation for Training Verifiably Robust Models, August 2019a. arXiv:1810.12715 [cs, stat].
- Gowal, S., Dvijotham, K., Stanforth, R., Bunel, R., Qin, C., Uesato, J., Arandjelovic, R., Mann, T. A., and Kohli, P. Scalable Verified Training for Provably Robust Image Classification. In *International Conference on Computer Vision*, 2019b.
- Jovanovic, N., Balunovic, M., Baader, M., and Vechev, M. T. On the Paradox of Certified Training. February 2021.
- Katz, G., Barrett, C. W., Dill, D. L., Julian, K., and Kochenderfer, M. J. Reluplex: An Efficient SMT Solver for Verifying Deep Neural Networks. In *Proceedings of the International Conference of Computer Aided Verification*, 2017.
- Kingma, D. P. and Ba, J. Adam: A Method for Stochastic Optimization. In *Proceedings of the International Conference on Learning Representations*, 2015.
- Krizhevsky, A. Learning multiple layers of features from tiny images. *citeSeer*, 2009.
- Kumar, A., Levine, A., Goldstein, T., and Feizi, S. Curse of Dimensionality on Randomized Smoothing for Certifiable Robustness. In *Proceedings of the International Conference on Machine Learning*, 2020.
- Le, Y. and Yang, X. S. Tiny imagenet visual recognition challenge. 2015.
- LeCun, Y., Bottou, L., Bengio, Y., and Haffner, P. Gradient-based learning applied to document recognition. *Proc. IEEE*, (11), 1998.
- Li, B., Chen, C., Wang, W., and Carin, L. Certified Adversarial Robustness with Additive Noise. In *Advances in Neural Information Processing Systems*, 2019.
- Lécuyer, M., Atlidakis, V., Geambasu, R., Hsu, D., and Jana, S. Certified Robustness to Adversarial Examples with Differential Privacy. In *IEEE Symposium on Security and Privacy*, 2019.

- 495 Madry, A., Makelov, A., Schmidt, L., Tsipras, D., and
 496 Vladu, A. Towards Deep Learning Models Resistant to
 497 Adversarial Attacks. In *Proceedings of the International*
 498 *Conference on Learning Representations*, 2018.
- 499
- 500 Mirman, M., Gehr, T., and Vechev, M. T. Differentiable
 501 Abstract Interpretation for Provably Robust Neural
 502 Networks. In *Proceedings of the International Conference*
 503 *on Machine Learning*, 2018.
- 504
- 505 Mohapatra, J., Ko, C.-Y., Weng, L., Chen, P.-Y., Liu, S., and
 506 Daniel, L. Hidden Cost of Randomized Smoothing. In
 507 *Proceedings of the International Conference on Artificial*
 508 *Intelligence and Statistics*, 2021.
- 509
- 510 Nagarajan, V. and Kolter, J. Z. Uniform convergence may be
 511 unable to explain generalization in deep learning. In *Advances*
 512 *in Neural Information Processing Systems*, 2019.
- 513 Nandy, J., Saha, S., Hsu, W., Lee, M. L., and Zhu, X. X.
 514 Towards Bridging the gap between Empirical and Certified
 515 Robustness against Adversarial Examples, July 2022.
 516 arXiv:2102.05096 [cs].
- 517
- 518 Raghunathan, A., Steinhardt, J., and Liang, P. Certified
 519 Defenses against Adversarial Examples. In *Proceedings*
 520 *of the International Conference on Learning Representations*, 2018.
- 521
- 522 Raghunathan, A., Xie, S. M., Yang, F., Duchi, J. C., and
 523 Liang, P. Understanding and Mitigating the Tradeoff
 524 between Robustness and Accuracy. In *Proceedings of the*
 525 *International Conference on Machine Learning*, 2020.
- 526
- 527 Salman, H., Yang, G., Zhang, H., Hsieh, C.-J., and Zhang,
 528 P. A Convex Relaxation Barrier to Tight Robustness
 529 Verification of Neural Networks. In *Advances in Neural*
 530 *Information Processing Systems*, 2019.
- 531
- 532 Sanyal, A., Hu, Y., and Yang, F. How unfair is private learning?
 533 In *Proceedings of the Conference on Uncertainty in*
 534 *Artificial Intelligence*, 2022.
- 535
- 536 Singh, G., Gehr, T., Mirman, M., Püschel, M., and Vechev,
 537 M. Fast and Effective Robustness Certification. In *Advances*
 538 *in Neural Information Processing Systems*, volume 31. Curran
 539 Associates, Inc., 2018.
- 540
- 541 Sun, J., Mehra, A., Kailkhura, B., Chen, P.-Y., Hendrycks,
 542 D., Hamm, J., and Mao, Z. M. A Spectral View of
 543 Randomized Smoothing Under Common Corruptions:
 544 Benchmarking and Improving Certified Robustness. In
Computer Vision – ECCV 2022, volume 13664. 2022.
- 545
- 546 Szegedy, C., Zaremba, W., Sutskever, I., Bruna, J., Erhan,
 547 D., Goodfellow, I. J., and Fergus, R. Intriguing properties
 548 of neural networks. In *Proceedings of the International*
 549 *Conference on Learning Representations*, 2014.
- Tsipras, D., Santurkar, S., Engstrom, L., Turner, A., and
 Madry, A. Robustness May Be at Odds with Accuracy. In
Proceedings of the International Conference on Learning
Representations, 2019.
- Weng, T., Zhang, H., Chen, H., Song, Z., Hsieh, C., Daniel,
 L., Boning, D. S., and Dhillon, I. S. Towards fast com-
 putation of certified robustness for relu networks. In
Proceedings of the International Conference on Machine
Learning, 2018.
- Wong, E. and Kolter, J. Z. Provable Defenses against Ad-
 versarial Examples via the Convex Outer Adversarial
 Polytope. In *Proceedings of the International Conference*
on Machine Learning, 2018.
- Wong, E., Schmidt, F. R., Metzen, J. H., and Kolter, J. Z.
 Scaling provable adversarial defenses. In *Advances in*
Neural Information Processing Systems, 2018.
- Xu, H., Liu, X., Li, Y., Jain, A. K., and Tang, J. To be
 Robust or to be Fair: Towards Fairness in Adversarial
 Training. In *Proceedings of the International Conference*
on Machine Learning, 2021.
- Xu, K., Shi, Z., Zhang, H., Wang, Y., Chang, K.-W., Huang,
 M., Kailkhura, B., Lin, X., and Hsieh, C.-J. Automatic
 Perturbation Analysis for Scalable Certified Robustness
 and Beyond. In *Advances in Neural Information Process-
 ing Systems*, 2020.
- Yang, G., Duan, T., Hu, J. E., Salman, H., Razenshteyn,
 I. P., and Li, J. Randomized Smoothing of All Shapes and
 Sizes. In *Proceedings of the International Conference on*
Machine Learning, pp. 10693–10705, 2020.
- Zhang, H., Weng, T.-W., Chen, P.-Y., Hsieh, C.-J., and
 Daniel, L. Efficient Neural Network Robustness Certifi-
 cation with General Activation Functions. In *Advances*
in Neural Information Processing Systems, 2018.
- Zhang, H., Yu, Y., Jiao, J., Xing, E. P., Ghaoui, L. E., and
 Jordan, M. I. Theoretically Principled Trade-off between
 Robustness and Accuracy. In *Proceedings of the Interna-
 tional Conference on Machine Learning*, 2019.
- Zhang, H., Chen, H., Xiao, C., Gowal, S., Stanforth, R.,
 Li, B., Boning, D. S., and Hsieh, C.-J. Towards Stable
 and Efficient Training of Verifiably Robust Neural Net-
 works. In *Proceedings of the International Conference*
on Learning Representations, 2020.

550 A. Certified defences for signal-directed perturbations

551 A.1. Certified defences for signal-directed perturbations

553 We now formulate the convex outer adversarial polytope (COAP) (Wong & Kolter, 2018) for adversaries that concentrate all
 554 their budget along the signal direction in the input. Our derivation can be seen as an extension of Wong & Kolter (2018);
 555 Erdemir et al. (2021).

557 **Network architecture** For the sake of clarity, we consider a 2-layers feed-forward ReLU network. However, our
 558 formulation can easily be extended to multiple layers. We define $f_\theta(x) : \mathbb{R}^d \rightarrow \mathbb{R}^2$ as follows:

$$559 \quad x \xrightarrow{x+\delta} z_1 \xrightarrow{W_1 z_1 + b_1} \hat{z}_2 \xrightarrow{\text{ReLU}(\cdot)} z_2 \xrightarrow{W_2 z_2 + b_2} \hat{z}_3 \quad (7)$$

560 where $x \in \mathbb{R}^d$, $z_1 \in \mathcal{B}_\epsilon(x)$, W_1 and W_2 are linear operators, and $\theta = \{W_i, b_i\}_{i=1,2}$ is the set of network parameters.

563 We define the adversarial polytope $\mathcal{Z}_\epsilon(x)$ as the set of all final-layer activations attainable by perturbing x with
 564 some $\tilde{x} \in \mathcal{B}_\epsilon(x)$:

$$565 \quad \mathcal{Z}_\epsilon(x) = \{f_\theta(\tilde{x}) : \tilde{x} \in \mathcal{B}_\epsilon(x)\}$$

567 The key idea behind COAP (Wong & Kolter, 2018) is to construct a convex outer bound to this adversarial polytope.
 568 Specifically, we relax the ReLU activations $z_2 = \text{ReLU}(\hat{z}_2)$ with their convex envelopes:

$$569 \quad z_2 \geq 0, \quad z_2 \geq \hat{z}_2, \quad (u - \ell)z_2 \leq u\hat{z}_2 - u\ell$$

570 where u and ℓ are respectively the pre-activations \hat{z}_2 upper and lower bounds. We assume that these bounds are known and
 571 provide a closed form solution to compute them in Appendix A.3. Further, let $\tilde{\mathcal{Z}}_\epsilon(x)$ be the outer bound to the adversarial
 572 polytope obtained from relaxing the ReLU constraints. Then, given a data point x with known label y , we can formalise the
 573 problem of finding an adversarial example with a linear program as follows:

$$574 \quad \min_{\hat{z}_3} [\hat{z}_3]_y - [\hat{z}_3]_{\bar{y}} = c^\top \hat{z}_3 \quad \text{s.t. } \hat{z}_3 \in \tilde{\mathcal{Z}}_\epsilon(x) \quad (8)$$

575 where \bar{y} is the binary negation of y . Note that if we solve this linear program and find that the objective is positive, then we
 576 know that no input perturbation within the threat model can misclassify the example. However, solving the linear program
 577 in Equation (8) for every example in the dataset is intractable. Therefore, we solve the dual formulation stated in the theorem
 578 below.

579 **Theorem 3.** *The dual of the linear program (8) can be written as*

$$580 \quad \begin{aligned} \max_{\alpha} \quad & \tilde{J}_\epsilon(x, g_\theta(c, \alpha)) \\ \text{s.t.} \quad & \alpha_j \in [0, 1], \forall j \end{aligned}$$

581 where $\tilde{J}_\epsilon(x, \nu_1, \nu_2, \nu_3)$ is equal to

$$582 \quad - \sum_{i=1}^2 \nu_{i+1}^\top b_i + \sum_{j \in \mathcal{I}} \ell_j [\nu_2]_j^+ - \hat{\nu}_1^\top x - \epsilon \|\hat{\nu}_1\|_1$$

583 and g_θ is a one-hidden layer neural network given by the equations

$$\begin{aligned} \nu_3 &= -c \\ \hat{\nu}_2 &= W_2^\top \nu_3 \\ [\nu_2]_j &= 0, \quad j \in \mathcal{I}^- \\ [\nu_2]_j &= [\hat{\nu}_2]_j, \quad j \in \mathcal{I}^+ \\ [\nu_2]_j &= \frac{u_j}{u_j - \ell_j} [\hat{\nu}_2]_j^+ - \alpha_j [\hat{\nu}_2]_j^-, \quad j \in \mathcal{I} \\ \hat{\nu}_1 &= W_1^\top \nu_2 \end{aligned}$$

584 where \mathcal{I}^- , \mathcal{I}^+ and \mathcal{I} denote the sets of activations in the hidden layer where ℓ and u are both negative, both positive or
 585 span zero, respectively.

586 Further, we take the feasible solution $\alpha_j = \frac{u_j}{u_j - \ell_j}$, as proposed in Wong & Kolter (2018). Hence, we can represent the dual
 587 problem as a linear back propagation network, which provides a tractable solution for a lower bound of the primal objective.

605 A.2. Proof of Theorem 3

 606 **Theorem 3.** *The dual of the linear program (8) can be written as*

608
 609
$$\max_{\alpha} \tilde{J}_\epsilon(x, g_\theta(c, \alpha))$$

 610
$$\text{s.t. } \alpha_j \in [0, 1], \forall j$$

 611 where $\tilde{J}_\epsilon(x, \nu_1, \nu_2, \nu_3)$ is equal to

612
 613
 614
$$-\sum_{i=1}^2 \nu_{i+1}^\top b_i + \sum_{j \in \mathcal{I}} \ell_j [\nu_2]_j^+ - \hat{\nu}_1^\top x - \epsilon \|\hat{\nu}_1\|_1$$

 615 and g_θ is a one-hidden layer neural network given by the equations

616
 617
 618
 619
 620
 621
 622
 623
 624
 625

$$\begin{aligned} \nu_3 &= -c \\ \hat{\nu}_2 &= W_2^\top \nu_3 \\ [\nu_2]_j &= 0, j \in \mathcal{I}^- \\ [\nu_2]_j &= [\hat{\nu}_2]_j, j \in \mathcal{I}^+ \\ [\nu_2]_j &= \frac{u_j}{u_j - \ell_j} [\hat{\nu}_2]_j^+ - \alpha_j [\hat{\nu}_2]_j^-, j \in \mathcal{I} \\ \hat{\nu}_1 &= W_1^\top \nu_2 \end{aligned}$$

 626 where \mathcal{I}^- , \mathcal{I}^+ and \mathcal{I} denote the sets of activations in the hidden layer where ℓ and u are both negative, both positive or span zero, respectively.

 627 *Proof.* Consider a data point x and let $\tilde{x} = x + \delta$ be the adversarial perturbed data point. First, we explicit all the constraints for the linear program defined in (8):

628
 629
 630
 631
 632

$$\begin{aligned} \min_{\hat{z}_3} [\hat{z}_3]_y - [\hat{z}_3]_{\bar{y}} &= c^\top \hat{z}_3, \quad \text{s.t.} \\ x + \delta &\in \mathcal{B}_\epsilon(x) \\ z_1 &= x + \delta \\ \hat{z}_2 &= W_1 z_1 + b_1 \\ \hat{z}_3 &= W_2 z_2 + b_2 \\ [z_2]_j &= 0, \forall j \in \mathcal{I}^- \\ [z_2]_j &= [\hat{z}_2]_j, \forall j \in \mathcal{I}^+ \\ [z_2]_j &\geq 0, \forall j \in \mathcal{I} \\ [z_2]_j &\geq [\hat{z}_2]_j, \forall j \in \mathcal{I} \\ ((u_j - \ell_j) [z_2]_j - u_j [\hat{z}_2]_j) &\leq -u_j \ell_j, \forall j \in \mathcal{I} \end{aligned}$$

 633 where \mathcal{I}^- , \mathcal{I}^+ and \mathcal{I} denote the sets of activations in the hidden layer where ℓ and u are both negative, both positive, or span zero respectively. In order to compute the dual of this problem, we associate the following Lagrangian variables with each of the constraints:

634
 635
 636
 637
 638
 639
 640
 641
 642
 643
 644
 645
 646

$$\begin{aligned} \hat{z}_2 &= W_1 z_1 + b_1 \Rightarrow \nu_2 \\ \hat{z}_3 &= W_2 z_2 + b_2 \Rightarrow \nu_3 \\ z_1 &= x + \delta \Rightarrow \psi \\ -[z_2]_j &\leq 0 \Rightarrow \mu_j, \forall j \in \mathcal{I} \\ [\hat{z}_2]_j - [z_2]_j &\leq 0 \Rightarrow \tau_j, \forall j \in \mathcal{I} \\ ((u_j - \ell_j) [z_2]_j - u_j [\hat{z}_2]_j) &\leq -u_j \ell_j \Rightarrow \lambda_j, \forall j \in \mathcal{I} \end{aligned}$$

 647 note that we do not define explicit dual variables for $[z_2]_j = 0$ and $[z_2]_j = [\hat{z}_2]_j$ as we can easily eliminate them. We write
 648
 649
 650
 651
 652
 653
 654
 655
 656
 657
 658
 659

660 the Lagrangian as follows:

$$\begin{aligned}
 663 \quad \mathcal{L}(z, \hat{z}, \nu, \delta, \lambda, \tau, \mu, \psi) = & - (W_1^\top \nu_2 + \psi)^\top z_1 - \sum_{j \in \mathcal{I}} \left(\mu_j + \tau_j - \lambda_j (u_j - \ell_j) + [W_2^\top \nu_3]_j \right) [z_2]_j \\
 664 \quad & + \sum_{j \in \mathcal{I}} (\tau_j - \lambda_j u_j + [\nu_2]_j) [\hat{z}_2]_j + (c + \nu_3)^\top \hat{z}_3 - \sum_{i=1}^2 \nu_{i+1}^\top b_i \\
 665 \quad & + \sum_{j \in \mathcal{I}} \lambda_j u_j \ell_j + \psi^\top x + \psi^\top \delta + \sum_{j \in \mathcal{I}^-} [\hat{z}_2]_j [\nu_2]_j \\
 666 \quad & + \sum_{j \in \mathcal{I}^+} [z_2]_j ([\nu_2]_j - [W_2^\top \nu_3]_j) \\
 667 \quad & \text{s.t. } \tilde{x} \in \mathcal{B}_\epsilon(x)
 \end{aligned}$$

673 and we take the infimum w.r.t. z, \hat{z}, δ :

$$\begin{aligned}
 678 \quad \inf_{z, \hat{z}, \delta} \mathcal{L}(z, \hat{z}, \nu, \delta, \lambda, \tau, \mu, \psi) = & - \inf_{z_2} \sum_{j \in \mathcal{I}} \left(\mu_j + \tau_j - \lambda_j (u_j - \ell_j) + [W_2^\top \nu_3]_j \right) [z_2]_j \\
 679 \quad & + \inf_{\hat{z}_2} \sum_{j \in \mathcal{I}} (\tau_j - \lambda_j u_j + [\nu_2]_j) [\hat{z}_2]_j + \inf_{\hat{z}_3} (c + \nu_3)^\top z_3 - \sum_{i=1}^2 \nu_{i+1}^\top b_i \\
 680 \quad & + \sum_{j \in \mathcal{I}} \lambda_j u_j \ell_j + \psi^\top x + \inf_{\tilde{x} \in \mathcal{B}_\epsilon(x)} \psi^\top \delta - \inf_{z_1} (W_1^\top \nu_2 + \psi)^\top z_1 \\
 681 \quad & + \inf_{\hat{z}_2} \sum_{j \in \mathcal{I}^-} [\hat{z}_2]_j [\nu_2]_j + \inf_{z_2} \sum_{j \in \mathcal{I}^+} [z_2]_j ([\nu_2]_j - [W_2^\top \nu_3]_j)
 \end{aligned}$$

689 Now, we can compute the infimum for the $\psi^\top \delta$ term:

$$692 \quad \inf_{\tilde{x} \in \mathcal{B}_\epsilon(x)} \psi^\top \delta = \inf_{\|\beta\|_1 \leq \epsilon} \psi_1 \cdot \beta = -\epsilon \cdot \|\psi_1\|_1$$

695 and since for all the other terms the infimum of a linear function is $-\infty$, except in the special case when it is identically
696 zero, the infimum of $\mathcal{L}(\cdot)$ becomes:

$$699 \quad \inf_{z, \hat{z}, \delta} \mathcal{L}(\cdot) = \begin{cases} -\sum_{i=1}^2 \nu_{i+1}^\top b_i + \sum_{j \in \mathcal{I}} \lambda_j u_j \ell_j + \psi^\top x - \epsilon \|\psi_1\|_1 & \text{if conditions} \\ -\infty & \text{else} \end{cases}$$

702 where the conditions to satisfy are:

$$\begin{aligned}
 705 \quad & \nu_3 = -c \\
 706 \quad & W_1^\top \nu_2 = -\psi \\
 707 \quad & [\nu_2]_j = 0, j \in \mathcal{I}_i^- \\
 708 \quad & [\nu_2]_j = [W_2^\top \nu_3]_j, j \in \mathcal{I}_i^+ \\
 709 \quad & \begin{aligned} (u_j - \ell_j) \lambda_j - \mu_j - \tau_j &= [W_2^\top \nu_3]_j \\ [\nu_2]_j &= u_j \lambda_j - \tau_j \end{aligned} \Big\} j \in \mathcal{I} \\
 710 \quad & \lambda, \tau, \mu \geq 0
 \end{aligned}$$

715 Thus, we can rewrite the dual problem as follows:

$$\begin{aligned}
 & \max_{\nu, \psi, \lambda, \tau, \mu} - \sum_{i=1}^2 \nu_{i+1}^\top b_i + \sum_{j \in \mathcal{I}} \lambda_j u_j \ell_j + \psi^\top x - \epsilon \|\psi_1\|_1 \\
 \text{s.t. } & \nu_3 = -c \\
 & W_1^\top \nu_2 = -\psi \\
 & [\nu_2]_j = 0, j \in \mathcal{I}_i^- \\
 & [\nu_2]_j = [W_2^\top \nu_3]_j, j \in \mathcal{I}_i^+ \\
 & \left. \begin{aligned} (u_j - \ell_j) \lambda_j - \mu_j - \tau_j &= [W_2^\top \nu_3]_j \\ [\nu_2]_j &= u_j \lambda_j - \tau_j \end{aligned} \right\} j \in \mathcal{I} \\
 & \lambda, \tau, \mu \geq 0
 \end{aligned}$$

730 Note that the dual variable λ corresponds to the upper bounds in the ReLU relaxation, while μ and τ correspond to the
 731 lower bounds. By the complementarity property, we know that at the optimal solution, these variables will be zero if the
 732 ReLU constraint is non-tight, or non-zero if the ReLU constraint is tight. Since the upper and lower bounds cannot be
 733 tight simultaneously, either λ or $\mu + \tau$ must be zero. This means that at the optimal solution to the dual problem we can
 734 decompose $[W_2^\top \nu_3]_j$ into positive and negative parts since $(u_j - \ell_j)\lambda_j \geq 0$ and $\tau_j + \mu_j \geq 0$:

$$\begin{aligned}
 (u_j - \ell_j)\lambda_j &= [W_2^\top \nu_3]_j^+ \\
 \tau_j + \mu_j &= [W_2^\top \nu_3]_j^-
 \end{aligned}$$

739 combining this with the constraint $[\nu_2]_j = u_j \lambda_j - \tau_j$ leads to

$$[\nu_2]_j = \frac{u_j}{u_j - \ell_j} [W_2^\top \nu_3]_j^+ - \alpha_j [W_2^\top \nu_3]_j^-$$

744 for $j \in \mathcal{I}$ and $0 \leq \alpha_j \leq 1$. Hence, we have that:

$$\lambda_j = \frac{u_j}{u_j - \ell_j} [\hat{\nu}_2]_j^+$$

749 Now, we denote $\hat{\nu}_1 = -\psi$ to make our notation consistent, and putting all of this together the dual objective becomes:

$$\begin{aligned}
 - \sum_{i=1}^2 \nu_{i+1}^\top b_i + \sum_{j \in \mathcal{I}} \lambda_j u_j \ell_j + \psi^\top x - \epsilon \|\psi_1\|_1 &= - \sum_{i=1}^2 \nu_{i+1}^\top b_i + \sum_{j \in \mathcal{I}} \frac{u_j \ell_j}{u_j - \ell_j} [\hat{\nu}_2]_j^+ - \hat{\nu}_1^\top x - \epsilon \|[\hat{\nu}_1]_1\|_1 \\
 &= - \sum_{i=1}^2 \nu_{i+1}^\top b_i + \sum_{j \in \mathcal{I}} \ell_j [\nu_2]_j^+ - \hat{\nu}_1^\top x - \epsilon \|[\hat{\nu}_1]_1\|_1
 \end{aligned}$$

757 and the final dual problem:

$$\begin{aligned}
 \max_{\nu, \hat{\nu}} & - \sum_{i=1}^2 \nu_{i+1}^\top b_i + \sum_{j \in \mathcal{I}} \ell_j [\nu_2]_j^+ - \hat{\nu}_1^\top x - \epsilon \|[\hat{\nu}_1]_1\|_1 \\
 \text{s.t. } & \nu_3 = -c \\
 & \hat{\nu}_2 = W_2^\top \nu_3 \\
 & [\nu_2]_j = 0, j \in \mathcal{I}_i^- \\
 & [\nu_2]_j = [\hat{\nu}_2]_j, j \in \mathcal{I}_i^+ \\
 & [\nu_2]_j = \frac{u_j}{u_j - \ell_j} [\hat{\nu}_2]_j^+ - \alpha_j [\hat{\nu}_2]_j^-, j \in \mathcal{I} \\
 & \hat{\nu}_1 = W_1^\top \nu_2
 \end{aligned}$$

□

770 **A.3. Computing upper and lower bounds on the pre-activations**

771 We address here the problem of obtaining the upper and lower bounds u and ℓ for the pre-activations \hat{z} . Specifically, the
772 following proposition gives a closed form solution.

773 **Proposition 4.** Consider the neural network f_θ defined in Equation (7). Let w_1 be the first column of W_1 . Then, for a data
774 point x and perturbation budget ϵ , we have the following element-wise bounds on the pre-activation vector \hat{z}_2 :

$$775 \quad \ell \leq \hat{z}_2 \leq u$$

776 where
777

$$778 \quad \ell = W_1x + b_1 - \epsilon|w_1| \text{ and } u = W_1x + b_1 + \epsilon|w_1|$$

781 *Proof.* Given a data point x and perturbation budget ϵ , let $\tilde{x} = x + \delta$ be the perturbed input to the network. First, we find an
782 upper bound the pre-activations values \hat{z}_2 :

$$783 \quad \hat{z}_2 = W_1(x + \delta) + b_1 = W_1x + b_1 + W_1\delta$$

784 In particular, we want to solve the following optimisation problem for each component of the pre-activation vector:
785

$$786 \quad u_i = \max_{\tilde{x} \in \mathcal{B}_\epsilon(x)} [\hat{z}_2]_i = [W_1x]_i + [b_1]_i + \max_{\tilde{x} \in \mathcal{B}_\epsilon(x)} [W_1\delta]_i$$

787 where u will be the vector containing element-wise upper bounds. Note that $\delta = \beta e_1$, thus the optimisation problem can be
788 rewritten as:

$$789 \quad \max_{\tilde{x} \in \mathcal{B}_\epsilon(x)} [W_1\delta]_i = \max_{\|\beta\|_1 \leq \epsilon} \beta \cdot [w_1]_i = \epsilon \cdot \| [w_1]_i \|_1$$

790 where w_1 is the first column of W_1 . The vector of upper bounds will then be:
791

$$792 \quad u = W_1x + b_1 + \epsilon|w_1|$$

793 Along the same lines, we can derive the vector of lower bounds ℓ :
794

$$795 \quad \ell = W_1x + b_1 - \epsilon|w_1|$$

796 □

801 **B. Theoretical results for signal-directed perturbations**

802 In this section, we prove that convex relaxations along the signal direction hurt robust generalisation. We focus on the
803 classification setting with binary cross-entropy loss:
804

$$805 \quad L(x, y) = y \log(\sigma(x)) + (1 - y) \log(1 - \sigma(x)) \quad (9)$$

806 where $\sigma(\cdot)$ is the sigmoid function.
807

808 First, in Lemma 1 we relate the robust error of the classifier f_θ to the signal parameter θ_1 . Specifically, we show
809 that robust error monotonically decreases in θ_1 .

810 **Lemma 1.** Let f_θ be the neural network defined in Equation (5) and \mathcal{B}_ϵ the threat model defined in Equation (4). We define
811 the robust risk \mathbf{R}_ϵ of f_θ as follows:
812

$$813 \quad \mathbf{R}_\epsilon(\theta) := \mathbb{P}_{(x,y)} [\exists z \in \mathcal{B}_\epsilon(x) : y \neq \text{sgn}(f_\theta(z))]$$

814 Then, $\mathbf{R}_\epsilon(\theta)$ is monotonically decreasing in θ_1 .
815

816 *Proof.* We assume, without loss of generality, that $\theta_1 > 0$, and since a and b are not trainable parameters we must have
817 $a > 0$ and $b < 0$ to solve the classification problem. The robust risk is then equal to:
818

$$819 \quad \begin{aligned} \mathbf{R}_\epsilon(\theta) &:= \mathbb{P}_{(x,y)} [\exists z \in \mathcal{B}_\epsilon(x) : y \neq \text{sgn}(f_\theta(z))] \\ &= \frac{1}{2} (\mathbb{P}_x [a \text{ReLU}(\theta^\top x) + b < 0 \mid y = 1] + \mathbb{P}_x [a \text{ReLU}(\theta^\top x) + b > 0 \mid y = -1]) \end{aligned}$$

further, we can remove the $\text{ReLU}(\cdot)$ using the fact that $\frac{b}{a} < 0$:

$$\begin{aligned}\mathbf{R}_\epsilon(\theta) &= \frac{1}{2} \left(\mathbb{P}_x \left[\{x : \theta^\top x + \frac{b}{a} < 0 \vee \theta^\top x < 0\} \mid y = 1 \right] + \mathbb{P}_x \left[\{x : \theta^\top x + \frac{b}{a} > 0 \wedge \theta^\top x > 0\} \mid y = -1 \right] \right) \\ &= \frac{1}{2} \left(\mathbb{P}_x \left[\theta^\top x + \frac{b}{a} < 0 \mid y = 1 \right] + \mathbb{P}_x \left[\theta^\top x + \frac{b}{a} > 0 \mid y = -1 \right] \right) \\ &= \frac{1}{2} \left(\mathbb{P}_x \left[\sum_{i=2}^d x_i \theta_i < -\theta_1(\gamma - \epsilon) - \frac{b}{a} \right] + \mathbb{P}_x \left[\sum_{i=2}^d x_i \theta_i > \theta_1(\gamma - \epsilon) - \frac{b}{a} \right] \right)\end{aligned}$$

Recall now that for the linearly separable distribution we have $\sum_{i=2}^d x_i \theta_i \sim \mathcal{N}(0, \sigma^2 \|\theta_{2:d}\|^2)$. Therefore, we can replace the probabilities with the standard normal CDF Φ :

$$\mathbf{R}_\epsilon(\theta) = \frac{1}{2} \left(\Phi \left(-\frac{(\gamma - \epsilon)\theta_1}{\sigma \|\theta_{2:d}\|_2} - \frac{b}{a\sigma \|\theta_{2:d}\|_2} \right) + \Phi \left(-\frac{(\gamma - \epsilon)\theta_1}{\sigma \|\theta_{2:d}\|_2} + \frac{b}{a\sigma \|\theta_{2:d}\|_2} \right) \right)$$

hence $\mathbf{R}_\epsilon(\theta)$ is monotonically decreasing in θ_1 and the statement follows. \square

B.1. Adversarial training

The basic idea behind adversarial training is to update the network parameters according to the following rule:

$$\theta \leftarrow \theta - \frac{\eta}{|D|} \sum_{(x,y) \in D} \nabla_\theta \max_{x+\delta \in \mathcal{B}_\epsilon(x)} L(f_\theta(x+\delta), y)$$

This is usually done by applying some first-order approximation to the maximisation problem. However, for our simplified network we can analytically compute the gradient. First of all, note that when L is the binary cross-entropy loss function we can rewrite the maximisation problem as follows:

$$\max_{x+\delta \in \mathcal{B}_\epsilon(x)} L(f_\theta(x+\delta), y) = L \left(\text{sgn}(y) \underbrace{\min_{x+\delta \in \mathcal{B}_\epsilon(x)} \text{sgn}(y) f_\theta(x+\delta)}_{:= J_\epsilon(x,y)}, y \right) \quad (10)$$

In particular, if J_ϵ is strictly positive then no adversarial example exists that fools the network. Further, note that this formulation is closely related to the objective considered in Appendix A.2 for the convex outer adversarial polytope. This will be useful when comparing COAP and AT gradients. Below we provide the gradient of the adversarial training objective w.r.t. the network parameters θ .

Proposition 5. Consider the neural network f_θ defined in Equation (5) and the threat model \mathcal{B}_ϵ defined in Equation (4). Let L be the binary cross-entropy loss function, as defined in Equation (9). Then, we have:

$$\begin{aligned}\nabla_{\theta_1} \max_{x+\delta \in \mathcal{B}_\epsilon(x)} L(f_\theta(x+\delta), y) \\ = -\text{sgn}(y) \boldsymbol{\sigma}(-J_\epsilon(x, y)) \begin{cases} a(x_1 - \epsilon \text{sgn}(\theta_1)) \mathbf{1}\{\ell > 0\} & \text{if } a \text{sgn}(y) > 0 \\ a(x_1 + \epsilon \text{sgn}(\theta_1)) \mathbf{1}\{u > 0\} & \text{if } a \text{sgn}(y) < 0 \end{cases}\end{aligned}$$

where $\ell = \theta^\top x - \epsilon |\theta_1|$ and $u = \theta^\top x + \epsilon |\theta_1|$ are respectively lower and upper bounds on the ReLU inputs.

Proof. Given a data point x with known label $y \in \{-1, 1\}$, when L is the binary cross-entropy loss function we have:

$$\max_{x+\delta \in \mathcal{B}_\epsilon(x)} L(f_\theta(x+\delta), y) = L \left(\text{sgn}(y) \min_{x+\delta \in \mathcal{B}_\epsilon(x)} \text{sgn}(y) f_\theta(x+\delta), y \right)$$

880 For our simplified network we can analytically compute a closed form solution of the minimisation problem:

$$\begin{aligned} J_\epsilon &:= \min_{x+\delta \in \mathcal{B}_\epsilon(x)} \operatorname{sgn}(y) (b + a \operatorname{ReLU}(\theta^\top(x + \delta))) \\ &= \begin{cases} \operatorname{sgn}(y) (b + a \max(0, \ell)) & \text{if } a \operatorname{sgn}(y) > 0 \\ \operatorname{sgn}(y) (b + a \max(0, u)) & \text{if } a \operatorname{sgn}(y) < 0 \end{cases} \\ &= \begin{cases} \operatorname{sgn}(y) (b + a \max(0, \ell)) & \text{if } a \operatorname{sgn}(y) > 0 \\ \operatorname{sgn}(y) (b + a \max(0, u)) & \text{if } a \operatorname{sgn}(y) < 0 \end{cases} \end{aligned}$$

891 where $\ell = \theta^\top x - \epsilon |\theta_1|$ and $u = \theta^\top x + \epsilon |\theta_1|$ are respectively lower and upper bounds on the pre-activations. Thus, we can
892 compute the gradients for adversarial training w.r.t the signal parameter:

$$\frac{\partial}{\partial \theta_1} J_\epsilon = \begin{cases} \operatorname{sgn}(y) a(x_1 - \epsilon \operatorname{sgn}(\theta_1)) \mathbf{1}\{\ell > 0\} & \text{if } a \operatorname{sgn}(y) > 0 \\ \operatorname{sgn}(y) a(x_1 + \epsilon \operatorname{sgn}(\theta_1)) \mathbf{1}\{u > 0\} & \text{if } a \operatorname{sgn}(y) < 0 \end{cases}$$

893 and applying the chain-rule we have:

$$\begin{aligned} \frac{\partial}{\partial \theta_1} L(\operatorname{sgn}(y) J_\epsilon, y) &= \frac{\partial}{\partial J_\epsilon} L(\operatorname{sgn}(y) J_\epsilon, y) \cdot \frac{\partial}{\partial \theta_1} J_\epsilon \\ &= \operatorname{sgn}(y) [\sigma(\operatorname{sgn}(y) J_\epsilon) - \mathbf{1}\{y = 1\}] \cdot \frac{\partial}{\partial \theta_1} J_\epsilon \\ &= -\operatorname{sgn}(y) \sigma(-J_\epsilon) \begin{cases} a(x_1 - \epsilon \operatorname{sgn}(\theta_1)) \mathbf{1}\{\ell > 0\} & \text{if } a \operatorname{sgn}(y) > 0 \\ a(x_1 + \epsilon \operatorname{sgn}(\theta_1)) \mathbf{1}\{u > 0\} & \text{if } a \operatorname{sgn}(y) < 0 \end{cases} \end{aligned}$$

907 where in the last equality we use a known property of the sigmoid function, $\sigma(x) = 1 - \sigma(-x)$. \square

B.2. Convex outer adversarial polytope

911 We now consider the dual approximation \tilde{J}_ϵ to the optimisation problem in Equation (10). Note that, for a binary classification
912 problem, we have $c = \operatorname{sgn}(y)$ and the dual objective in Theorem 3 becomes:

$$\tilde{J}_\epsilon(x, g_\theta(c, \alpha)) = \tilde{J}_\epsilon(x, y) \quad (11)$$

915 where we set α to the dual feasible solution and for the sake of clarity we omit the dependence on the network parameters θ .

916 We are particularly interested in the data points for which $J_\epsilon(x, y) \neq \tilde{J}_\epsilon(x, y)$, i.e., when the certified and adversarial
917 training objectives differ. Below, we provide a necessary and sufficient condition to have a mismatch between the two
918 objectives.

919 **Proposition 6.** Consider the neural network f_θ defined in Equation (5) and the threat model \mathcal{B}_ϵ defined in Equation (4).
920 Let L be the binary cross-entropy loss function, as defined in Equation (9). Further, we define $\ell = \theta^\top x - \epsilon |\theta_1|$ and
921 $u = \theta^\top x + \epsilon |\theta_1|$ respectively as lower and upper bounds on the ReLU inputs. Let $\mathcal{I}^* = \{(x, y) : 0 \in [\ell, u] \wedge a \operatorname{sgn}(y) > 0\}$.
922 Then, for data points in \mathcal{I}^* , we have that AT and COAP gradients differ:

$$\nabla_{\theta_1} J_\epsilon(x, y) \neq \nabla_{\theta_1} \tilde{J}_\epsilon(x, y) \quad \forall (x, y) \in \mathcal{I}^*$$

923 and COAP gradient is given by:

$$\begin{aligned} \nabla_{\theta_1} L(\operatorname{sgn}(y) \tilde{J}_\epsilon(x, y), y) &= -\frac{a \operatorname{sgn}(y) \sigma(-\tilde{J}_\epsilon(x, y))}{2\epsilon} \left(\frac{\ell}{\|\theta_1\|_1} (x_1 + \epsilon \operatorname{sgn}(\theta_1)) + u \frac{x_1 \|\theta_1\|_1 - \theta^\top x \operatorname{sgn}(\theta_1)}{\theta_1^2} \right) \end{aligned}$$

932 Further, for data points that are not in \mathcal{I}^* we have that AT and COAP gradients are equivalent:

$$\nabla_{\theta_1} J_\epsilon(x, y) = \nabla_{\theta_1} \tilde{J}_\epsilon(x, y) \quad \forall (x, y) \notin \mathcal{I}^*$$

935 *Proof.* For the sake of clarity, we report here the definition of COAP objective from Appendix A.2.

$$936 \quad 937 \quad 938 \quad 939 \quad \tilde{J}_\epsilon(x, y) = -\sum_{i=1}^2 \nu_{i+1}^\top b_i + \sum_{j \in \mathcal{I}} \ell_j [\nu_2]_j^+ - \hat{\nu}_1^\top x - \epsilon \|\hat{\nu}_1\|_1$$

940 Further, recall that the dual variables ν are given by the following equations:

$$941 \quad 942 \quad \nu_3 = -c \\ 943 \quad \hat{\nu}_2 = W_2^\top \nu_3 \\ 944 \quad [\nu_2]_j = 0, \quad j \in \mathcal{I}^- \\ 945 \quad [\nu_2]_j = [\hat{\nu}_2]_j, \quad j \in \mathcal{I}^+ \\ 946 \quad [\nu_2]_j = \frac{u_j}{u_j - \ell_j} [\hat{\nu}_2]_j^+ - \alpha_j [\hat{\nu}_2]_j^-, \quad j \in \mathcal{I} \\ 947 \quad \hat{\nu}_1 = W_1^\top \nu_2$$

948 where \mathcal{I}^- , \mathcal{I}^+ and \mathcal{I} denote the sets of activations in the hidden layer where ℓ and u are both negative, both positive and span zero, respectively.

949 First, we consider the case when the neuron is always dead, i.e., $\ell < u < 0$. The dual variables are:

$$950 \quad 951 \quad \nu_3 = -\text{sgn}(y) \\ 952 \quad \hat{\nu}_2 = -a \text{sgn}(y) \\ 953 \quad \nu_2 = 0 \\ 954 \quad \hat{\nu}_1 = 0$$

955 Hence, AT and COAP objectives are equal in this case:

$$956 \quad 957 \quad \tilde{J}_\epsilon = \text{sgn}(y)b = J_\epsilon$$

958 where the last equality follows from Appendix B.1.

959 Next, we consider the case when the neuron is always active, i.e., $0 < \ell < u$. The dual variables are:

$$960 \quad 961 \quad \nu_3 = -\text{sgn}(y) \\ 962 \quad \hat{\nu}_2 = -a \text{sgn}(y) \\ 963 \quad \nu_2 = -a \text{sgn}(y) \\ 964 \quad \hat{\nu}_1 = -a \text{sgn}(y) \cdot \theta$$

965 and the dual objective becomes:

$$966 \quad 967 \quad \tilde{J}_\epsilon = -\nu_3^\top b - \hat{\nu}_1^\top x - \epsilon \|\hat{\nu}_1\|_1 \\ 968 \quad = \text{sgn}(y) (b + a(\theta^\top x)) - \epsilon \|a \text{sgn}(y)\theta\|_1 \\ 969 \quad = \begin{cases} \text{sgn}(y) (b + a\ell) & \text{if } a \text{sgn}(y) > 0 \\ \text{sgn}(y) (b + au) & \text{if } a \text{sgn}(y) < 0 \end{cases} \\ 970 \quad = J_\epsilon$$

971 where the last equality follows from the fact that $0 < \ell < u$.

972 Finally, we consider the case when the neuron is in the activation set \mathcal{I} , i.e., $\ell < 0 < u$. The dual variables are:

$$973 \quad 974 \quad \nu_3 = -\text{sgn}(y) \\ 975 \quad \hat{\nu}_2 = -a \text{sgn}(y) \\ 976 \quad \nu_2 = -a \text{sgn}(y) \frac{u}{2\epsilon \|\theta\|_1} \\ 977 \quad \hat{\nu}_1 = -a \text{sgn}(y) \frac{u}{2\epsilon \|\theta\|_1} \cdot \theta$$

990 Here we have two cases, when $\hat{\nu}_2 > 0$ we can rewrite the dual objective as:

$$991 \quad 992 \quad 993 \quad 994 \quad 995 \quad 996 \quad 997 \quad 998 \quad 999 \quad 1000 \quad 1001 \quad 1002 \quad 1003 \quad 1004 \quad 1005 \quad 1006 \quad 1007 \quad 1008 \quad 1009 \quad 1010 \quad 1011 \quad 1012 \quad 1013 \quad 1014 \quad 1015 \quad 1016 \quad 1017 \quad 1018 \quad 1019 \quad 1020 \quad 1021 \quad 1022 \quad 1023 \quad 1024 \quad 1025 \quad 1026 \quad 1027 \quad 1028 \quad 1029 \quad 1030 \quad 1031 \quad 1032 \quad 1033 \quad 1034 \quad 1035 \quad 1036 \quad 1037 \quad 1038 \quad 1039 \quad 1040 \quad 1041 \quad 1042 \quad 1043 \quad 1044 \quad \widetilde{J}_\epsilon = \text{sgn}(y)(b + au) = J_\epsilon$$

and the two objectives coincide.

When $\nu_2 < 0$ we can rewrite the dual objective as:

$$\widetilde{J}_\epsilon = \text{sgn}(y) \left(b + \frac{au\ell}{2\epsilon \|\theta_1\|_1} \right) \neq J_\epsilon$$

Hence, the only case when COAP gradient differs from AT gradient is when $\nu_2 < 0$ and the neuron belongs to the activation set \mathcal{I} .

We compute the partial derivative w.r.t. the signal parameter θ_1 in this case, by the chain rule we have:

$$\begin{aligned} \frac{\partial}{\partial \theta_1} L(\text{sgn}(y) \cdot \widetilde{J}_\epsilon, y) \\ = \frac{\partial}{\partial \widetilde{J}_\epsilon} L(\text{sgn}(y) \cdot \widetilde{J}_\epsilon, y) \cdot \frac{\partial}{\partial \theta_1} \widetilde{J}_\epsilon \\ = \text{sgn}(y) [\sigma(\text{sgn}(y) \cdot \widetilde{J}_\epsilon) - \mathbf{1}\{y=1\}] \cdot \frac{\partial}{\partial \theta_1} \widetilde{J}_\epsilon \\ = -\frac{a \text{sgn}(y) \sigma(-\widetilde{J}_\epsilon)}{2\epsilon} \left(\frac{\ell}{\|\theta_1\|_1} (x_1 + \epsilon \text{sgn}(\theta_1)) + u \frac{x_1 \|\theta_1\|_1 - \theta^\top x \text{sgn}(\theta_1)}{\theta_1^2} \right) \end{aligned}$$

□

B.3. Proof of Theorem 1

Theorem 1. Let \mathbf{R}_ϵ be the robust risk of f_θ , defined as

$$1019 \quad \mathbf{R}_\epsilon(\theta) := \mathbb{P}_{(x,y)} [\exists z \in \mathcal{B}_\epsilon(x) : y \neq \text{sgn}(f_\theta(z))].$$

Let $\bar{\theta}$ and $\tilde{\theta}$ be the network parameters after one step of gradient descent with respect to AT and COAP objectives. If the initialised network parameters θ satisfy

$$1024 \quad 1025 \quad 1026 \quad 1027 \quad \frac{\|\theta_{2:d}\|_2}{\|\theta_1\|_2} > \sqrt{\max \left(\frac{(7\epsilon - \gamma)(\gamma + \epsilon)^4}{4\sigma^2(\gamma^2 - 10\gamma\epsilon + 13\epsilon^2)}, \frac{(\gamma + \epsilon)^3}{12\sigma^2\epsilon} \right)}$$

and $\frac{2}{3}\gamma < \epsilon < \gamma$, COAP yields higher robust risk than AT:

$$1030 \quad 1031 \quad \mathbf{R}_\epsilon(\tilde{\theta}) > \mathbf{R}_\epsilon(\bar{\theta}).$$

Proof. First we assume, without loss of generality, that at initialisation $\theta_1 > 0$, and since a and b are not trainable parameters we must have $a > 0$ and $b < 0$ to include the ground truth in our hypothesis class.

Let J_ϵ be the adversarial training inner maximisation as defined in Equation (10). Then, AT solves the following optimisation problem:

$$1037 \quad 1038 \quad \min_{\theta} \mathbb{E}_{(x,y)} [L(\sigma(\text{sgn}(y)J_\epsilon(x,y)), y)]$$

Similarly, let \widetilde{J}_ϵ be the COAP dual approximation to the inner maximization described in Equation (11). Then, COAP solves the following optimisation problem:

$$1042 \quad 1043 \quad 1044 \quad \min_{\theta} \mathbb{E}_{(x,y)} [L(\sigma(\text{sgn}(y)\widetilde{J}_\epsilon), y)]$$

1045 Since we are only training the signal parameter θ_1 , after one gradient descent step, we have:

$$\|\bar{\theta}_{2:d}\|_2 = \|\tilde{\theta}_{2:d}\|_2$$

1049 Further, from Lemma 1 we know that AT yields smaller robust risk than COAP if the following holds:

$$\bar{\theta}_1 > \tilde{\theta}_1 \implies \mathbf{R}_\epsilon(\tilde{\theta}) > \mathbf{R}_\epsilon(\bar{\theta})$$

1052 which, after one step of gradient descent, is equivalent to:

$$\mathbb{E}_{(x,y)} [\nabla_{\theta_1} L (\boldsymbol{\sigma} (\text{sgn}(y) J_\epsilon(x, y)), y)] < \mathbb{E}_{(x,y)} [\nabla_{\theta_1} L (\boldsymbol{\sigma} (\text{sgn}(y) \tilde{J}_\epsilon(x, y)), y)]$$

1055 Now recall from Theorems 5 and 6 that the gradients of AT and COAP differ only on the set \mathcal{I}^* . In particular, we have that:

$$(x, y) \notin \mathcal{I}^* \implies \nabla_{\theta_1} L (\boldsymbol{\sigma} (\text{sgn}(y) J_\epsilon(x, y)), y) = \nabla_{\theta_1} L (\boldsymbol{\sigma} (\text{sgn}(y) \tilde{J}_\epsilon(x, y)), y) < 0$$

1059 and

$$(x, y) \in \mathcal{I}^* \implies 0 = \nabla_{\theta_1} L (\boldsymbol{\sigma} (\text{sgn}(y) J_\epsilon(x, y)), y) \neq \nabla_{\theta_1} L (\boldsymbol{\sigma} (\text{sgn}(y) \tilde{J}_\epsilon(x, y)), y)$$

1062 Hence, for our purpose we need to show that:

$$\mathbb{E}_{(x,y)} [\nabla_{\theta_1} L (\boldsymbol{\sigma} (\text{sgn}(y) \tilde{J}_\epsilon(x, y)), y) \mid (x, y) \in \mathcal{I}^*] > 0 \quad (12)$$

1066 Our strategy will be to lower-bound the expectation in Equation (12) with some strictly positive quantity. We define

$$Z = \sum_{i=2}^d \theta_i x_i$$

1071 and plug-in the gradient computed in Theorem 6:

$$\begin{aligned} & \mathbb{E}_{(x,y)} [\nabla_{\theta_1} L (\boldsymbol{\sigma} (\text{sgn}(y) \tilde{J}_\epsilon(x, y)), y) \mid (x, y) \in \mathcal{I}^*] \\ &= \mathbb{E}_{(x,y)} \left[\frac{a \boldsymbol{\sigma} (-\tilde{J}_\epsilon(x, y))}{2\epsilon} \left(-\frac{\ell}{\theta_1} (\gamma + \epsilon) + u \frac{\sum_{i=2}^d x_i \theta_i}{\theta_1^2} \right) \mid (x, y) \in \mathcal{I}^* \right] \\ &= \frac{a}{2\theta_1 \epsilon} \mathbb{E}_{(x,y)} [\boldsymbol{\sigma} (-\tilde{J}_\epsilon(x, y)) \left(-\ell(\gamma + \epsilon) + u \frac{Z}{\theta_1} \right) \mid (x, y) \in \mathcal{I}^*] \\ &= \frac{a}{2\theta_1 \epsilon} \mathbb{E}_{(x,y)} [\boldsymbol{\sigma} (-\tilde{J}_\epsilon(x, y)) u \frac{Z}{\theta_1} - \boldsymbol{\sigma} (-\tilde{J}_\epsilon(x, y)) \ell(\gamma + \epsilon) \mid (x, y) \in \mathcal{I}^*] \end{aligned}$$

1083 Now, we observe that Z is always negative on the set \mathcal{I}^* , since we need to satisfy the constraint $\ell < 0 < u$:

$$(x, y) \in \mathcal{I}^* \implies -\theta_1(\gamma + \epsilon) < \sum_{i=2}^d \theta_i x_i < -\theta_1(\gamma - \epsilon) < 0$$

1087 Further, from Appendix B.2 we have:

$$(x, y) \in \mathcal{I}^* \implies \boldsymbol{\sigma} (-\tilde{J}_\epsilon(x, y)) \geq \frac{1}{2}$$

1091 Combining these two observations we can lower-bound the expectation:

$$\begin{aligned} & \mathbb{E}_{(x,y)} [\nabla_{\theta_1} L (\boldsymbol{\sigma} (\text{sgn}(y) \tilde{J}_\epsilon(x, y)), y) \mid (x, y) \in \mathcal{I}^*] \\ &= \frac{a}{2\theta_1 \epsilon} \mathbb{E}_{(x,y)} [\boldsymbol{\sigma} (-\tilde{J}_\epsilon(x, y)) u \frac{Z}{\theta_1} - \boldsymbol{\sigma} (-\tilde{J}_\epsilon(x, y)) \ell(\gamma + \epsilon) \mid (x, y) \in \mathcal{I}^*] \\ &\geq \frac{a}{2\theta_1 \epsilon} \mathbb{E}_{(x,y)} \left[u \frac{Z}{\theta_1} - \frac{\gamma + \epsilon}{2} \ell \mid (x, y) \in \mathcal{I}^* \right] \end{aligned}$$

1100 Now, we need to show that this lower-bound is strictly positive:

$$1102 \quad \mathbb{E}_{(x,y)} \left[u \frac{Z}{\theta_1} - \frac{\gamma + \epsilon}{2} \ell \mid (x,y) \in \mathcal{I}^* \right] > 0$$

1104 Note that, we can further expand this expression:

$$1106 \quad \begin{aligned} \mathbb{E}_{(x,y)} \left[u \frac{Z}{\theta_1} - \frac{\gamma + \epsilon}{2} \ell \mid (x,y) \in \mathcal{I}^* \right] \\ 1109 \quad = -(\gamma^2 - \epsilon^2)\theta_1^2 + (\gamma + \epsilon)\theta_1 \mathbb{E}[Z \mid (x,y) \in \mathcal{I}^*] + 2\mathbb{E}[Z^2 \mid (x,y) \in \mathcal{I}^*] \end{aligned}$$

1111 Further, $Z \mid (x,y) \in \mathcal{I}^*$ is distributed as a truncated normal with:

$$1113 \quad \alpha = -\frac{\theta_1(\gamma + \epsilon)}{\sigma \|\theta_{2:d}\|_2} \text{ and } \beta = -\frac{\theta_1(\gamma - \epsilon)}{\sigma \|\theta_{2:d}\|_2}$$

1115 Hence, we can plug in the expectations of the truncated normal distribution to obtain the following:

$$\begin{aligned} 1119 \quad & -(\gamma^2 - \epsilon^2)\theta_1^2 + \theta_1(\gamma + \epsilon)\mathbb{E}[Z \mid (x,y) \in \mathcal{I}^*] + 2\mathbb{E}[Z^2 \mid (x,y) \in \mathcal{I}^*] \\ 1120 \quad & = -(\gamma^2 - \epsilon^2)\theta_1^2 + 2\sigma^2 \|\theta_{2:d}\|_2^2 + \sigma \|\theta_{2:d}\|_2 \theta_1 \frac{(\gamma - 3\epsilon)\phi(\beta) - (\gamma + \epsilon)\phi(\alpha)}{\Phi(\beta) - \Phi(\alpha)} \\ 1123 \quad & \propto -(\gamma^2 - \epsilon^2) + 2\sigma^2 r^2 + \sigma r \frac{(\gamma - 3\epsilon)\phi(\beta) - (\gamma + \epsilon)\phi(\alpha)}{\Phi(\beta) - \Phi(\alpha)} \\ 1125 \quad & = -f(r) \end{aligned}$$

1127 where we define $r = \frac{\|\theta_{2:d}\|_2}{\|\theta_1\|_2}$. Now, under our assumptions, from Lemma 4 we have:

$$1130 \quad f(r) < 0, \forall r > \sqrt{\max \left(\frac{(7\epsilon - \gamma)(\gamma + \epsilon)^4}{4\sigma^2(\gamma^2 - 10\gamma\epsilon + 13\epsilon^2)}, \frac{(\gamma + \epsilon)^3}{12\sigma^2\epsilon} \right)}$$

1133 and hence it follows that $\mathbf{R}_\epsilon(\tilde{\theta}) > \mathbf{R}_\epsilon(\bar{\theta})$.

1134 Next, we prove that the accuracy discrepancy is higher for COAP, that is:

$$1136 \quad \max_k \mathbf{R}^k(\tilde{\theta}) - \mathbf{R}(\tilde{\theta}) > \max_k \mathbf{R}^k(\bar{\theta}) - \mathbf{R}(\bar{\theta})$$

1138 As we already showed that $\bar{\theta}_1 > \tilde{\theta}_1$, our strategy will be to prove that the accuracy discrepancy is a decreasing function of θ_1 . First, note that for $\theta_1 >$ we have from Lemma 1:

$$\begin{aligned} 1142 \quad \max_k \mathbf{R}^k(\theta) - \mathbf{R}(\theta) &= \mathbf{R}_\epsilon^1(\theta) - \frac{\mathbf{R}_\epsilon^1(\theta) - \mathbf{R}_\epsilon^{-1}(\theta)}{2} \\ 1144 \quad &\propto \mathbf{R}_\epsilon^1(\theta) - \mathbf{R}_\epsilon^{-1}(\theta) := a(\theta) \end{aligned}$$

1145 Next, we take the derivative of $a(\theta)$ w.r.t. θ_1 and show that it is negative:

$$\begin{aligned} 1148 \quad \frac{\partial}{\partial \theta_1} a(\theta) &= \frac{\partial}{\partial \theta_1} [\Phi(-c_1\theta_1 + c_2) - \Phi(-c_1\theta_1 - c_2)] \\ 1149 \quad &= c_1 [\phi(-c_1\theta_1 - c_2) - \phi(-c_1\theta_1 + c_2)] \\ 1151 \quad &< 0 \quad \forall \theta_1, c_1, c_2 > 0 \end{aligned}$$

1153 which concludes the proof. □

1154

C. Auxiliary lemmas
C.1. Upper bound on the exponential function

Lemma 2. Let $f : \mathbb{R} \rightarrow \mathbb{R}$ be the function defined by $f(x) = \exp(x)$. When $x \leq 0$ and n is even we have:

$$f(x) \leq 1 + x + \frac{x^2}{2!} + \cdots + \frac{x^n}{n!}$$

Proof. Let $g : (-\infty, 0] \rightarrow \mathbb{R}$ be the function defined by

$$g(x) = 1 + x + \frac{x^2}{2!} + \cdots + \frac{x^n}{n!} - \exp(x)$$

Since $g(x) \rightarrow \infty$ as $x \rightarrow -\infty$, g must attain an absolute minimum somewhere on the interval $(-\infty, 0]$. Now, differentiating we have:

- If g has an absolute minimum at 0, then for all x , $g(x) \geq g(0) = 1 - \exp(0) = 0$, so we are done.

- If g has an absolute minimum at y for some $y < 0$, then $g'(y) = 0$. But differentiating,

$$g'(y) = 1 + y + \frac{y^2}{2!} + \cdots + \frac{y^{n-1}}{(n-1)!} - \exp(y) = g(y) - \frac{y^n}{n!}.$$

Therefore, for any x ,

$$g(x) \geq g(y) = \frac{y^n}{n!} + g'(y) = \frac{y^n}{n!} > 0,$$

since n is even.

□

C.2. Lower bound on the difference of Gaussian CDFs

Lemma 3. Let $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ be the function defined by $f(x, y) = \Phi(y) - \Phi(x)$. When $x < y < 0$ we have:

$$\phi(0) \left(y - x + \frac{x^3}{6} \right) \leq \Phi(y) - \Phi(x)$$

where Φ and ϕ are respectively the CDF and PDF of the standard Gaussian distribution.

Proof. First, we want to prove that $\frac{2x}{\sqrt{\pi}}$ is a lower bound for the error function $\text{erf}(x)$ when $x \leq 0$. That is, we want to show that $f(x) \geq 0$ where $f : (-\infty, 0] \rightarrow \mathbb{R}$ is the function defined by:

$$f(x) = \text{erf}(x) - \frac{2x}{\sqrt{\pi}}$$

Since f is continuous and $f(x) \rightarrow \infty$ as $x \rightarrow -\infty$, f must attain an absolute minimum on the interval $(-\infty, 0]$. Now, differentiating we have:

$$f'(x) = \frac{2}{\sqrt{\pi}} \exp(-x^2) - \frac{2}{\sqrt{\pi}}$$

hence f attains an absolute minimum at 0 and we have $f(x) \geq f(0) = 0$.

Next, we show that $\frac{2}{\sqrt{\pi}}(x - x^3/3)$ is an upper bound for $\text{erf}(x)$ when $x \leq 0$. Let $g : (-\infty, 0] \rightarrow \mathbb{R}$ the function defined by:

$$g(x) = \frac{2}{\sqrt{\pi}}(x - x^3/3) - \text{erf}(x)$$

Similarly, since g is continuous and $g(x) \rightarrow \infty$ as $x \rightarrow -\infty$, g must attain an absolute minimum on the interval $(-\infty, 0]$. Now, differentiating we have:

$$g'(x) = \frac{2}{\sqrt{\pi}}(1 - x^2 - \exp(-x^2))$$

1210 hence g attains an absolute minimum at 0 and we have $g(x) \geq g(0) = 0$.

1211 Now, since $a < b < 0$ we can use the erf bounds derived above:

$$\begin{aligned} 1213 \quad \Phi(b) - \Phi(a) &= \frac{1}{2} \left(\operatorname{erf}(b/\sqrt{2}) - \operatorname{erf}(a/\sqrt{2}) \right) \\ 1214 \quad &\geq \frac{1}{\sqrt{\pi}} \left(\frac{b}{\sqrt{2}} - \frac{a}{\sqrt{2}} + \frac{a^3}{6\sqrt{2}} \right) \\ 1215 \quad &= \phi(0) \left(b - a + \frac{a^3}{6} \right) \\ 1216 \quad & \\ 1217 \quad & \\ 1218 \quad & \\ 1219 \end{aligned}$$

1220 which concludes the proof. \square

1221

1222 C.3. Upper bound on the ratio of Gaussian PDFs and CDFs

1223

1224 **Lemma 4.** Suppose $f : \mathbb{R} \rightarrow \mathbb{R}$ is defined as follows:

1225

$$1226 \quad f(r) = \gamma^2 - \epsilon^2 - 2\sigma^2 r^2 - \sigma r \frac{(\gamma - 3\epsilon)\phi(\beta) - (\gamma + \epsilon)\phi(\alpha)}{\Phi(\beta) - \Phi(\alpha)}$$

1227

1228 where $\alpha := -\frac{\gamma+\epsilon}{r\sigma}$, $\beta := -\frac{\gamma-\epsilon}{r\sigma}$, Φ and ϕ are respectively the standard Gaussian CDF and PDF.

1229

1230 Assume that:

1231

$$1232 \quad \frac{5+2\sqrt{3}}{13} \gamma < \epsilon < \gamma$$

1233

1234 Then, we have:

1235

$$1236 \quad f(r) < 0, \quad \forall r > \sqrt{\max \left(\frac{(7\epsilon - \gamma)(\gamma + \epsilon)^4}{4\sigma^2(\gamma^2 - 10\gamma\epsilon + 13\epsilon^2)}, \frac{(\gamma + \epsilon)^3}{12\sigma^2\epsilon} \right)}$$

1237

1238 *Proof.* We begin by providing a lower bound on the difference of gaussian cdfs. Applying Lemma 3 with $x = \alpha$ and $y = \beta$ we have:

1239

$$1240 \quad \Phi(\beta) - \Phi(\alpha) \geq \left(\frac{2\epsilon}{r\sigma} - \frac{(\gamma + \epsilon)^3}{6\sigma^3 r^3} \right) \phi(0), \quad \alpha < \beta < 0$$

1241

1242

1243

1244

1245 Next, we can upper-bound f :

1246

$$\begin{aligned} 1247 \quad f(r) &\leq \gamma^2 - \epsilon^2 - 2\sigma^2 r^2 - \sigma r \frac{(\gamma - 3\epsilon)\phi(\beta) - (\gamma + \epsilon)\phi(\alpha)}{\left(\frac{2\epsilon}{r\sigma} - \frac{(\gamma + \epsilon)^3}{6\sigma^3 r^3} \right) \phi(0)} \\ 1248 \quad &\leq \gamma^2 - \epsilon^2 - 2\sigma^2 r^2 - \sigma^2 r^2 \frac{(\gamma - 3\epsilon)\phi(0) - (\gamma + \epsilon)\phi(\alpha)}{\left(2\epsilon - \frac{(\gamma + \epsilon)^3}{6r^2\sigma^2} \right) \phi(0)} \\ 1249 \quad &= \gamma^2 - \epsilon^2 - 2\sigma^2 r^2 - \sigma^2 r^2 \frac{(\gamma - 3\epsilon) - (\gamma + \epsilon) \exp(-\alpha^2/2)}{2\epsilon - \frac{(\gamma + \epsilon)^3}{6\sigma^2 r^2}} \\ 1250 \quad & \\ 1251 \quad & \\ 1252 \quad & \\ 1253 \quad & \\ 1254 \quad & \\ 1255 \quad & \\ 1256 \quad & \\ 1257 \quad & \\ 1258 \quad & \\ 1259 \quad & \\ 1260 \quad & \\ 1261 \quad & \\ 1262 \quad & \\ 1263 \quad & \\ 1264 \quad & \end{aligned}$$

Now, we use the upper-bound for the exponential function from Lemma 2 with $n = 2$:

$$\exp(x) \leq 1 + x - x^2/2, \quad \forall x \leq 0$$

and substituting it back into our upper-bound for f we get:

1261

1262

1263

1264

$$f(r) \leq \gamma^2 - \epsilon^2 - 2\sigma^2 r^2 - \sigma^2 r^2 \frac{(\gamma - 3\epsilon) - (\gamma + \epsilon) \left(1 - \frac{(\gamma + \epsilon)^2}{2r^2\sigma^2} + \frac{(\gamma + \epsilon)^4}{8r^4\sigma^4} \right)}{2\epsilon - \frac{(\gamma + \epsilon)^3}{6r^2\sigma^2}}$$

1265 which can be further simplified:

$$\begin{aligned} f(r) &\leq \gamma^2 - \epsilon^2 - 2\sigma^2 r^2 - \sigma^2 r^2 \frac{(\gamma - 3\epsilon) - (\gamma + \epsilon)(1 - \frac{(\gamma+\epsilon)^2}{2r^2\sigma^2} + \frac{(\gamma+\epsilon)^4}{8r^4\sigma^4})}{2\epsilon - \frac{(\gamma+\epsilon)^3}{6r^2\sigma^2}} \\ &= \frac{(\gamma - 7\epsilon)(\gamma + \epsilon)^4 + 4r^2\sigma^2(\gamma + \epsilon)(\gamma^2 - 10\gamma\epsilon + 13\epsilon^2)}{4(\gamma + \epsilon)^3 - 48r^2\sigma^2\epsilon} \\ &= u(r) \end{aligned}$$

1273 and we have that for $\epsilon > \frac{5+2\sqrt{3}}{13}\gamma$ and $r > \sqrt{\max\left(\frac{(7\epsilon-\gamma)(\gamma+\epsilon)^4}{4\sigma^2(\gamma^2-10\gamma\epsilon+13\epsilon^2)}, \frac{(\gamma+\epsilon)^3}{12\sigma^2\epsilon}\right)}$ the upper bound is negative, i.e. $u(r) < 0$. \square

1277 D. Experimental details

1279 D.1. Synthetic experiments with signal-directed perturbations

1281 Below we provide detailed experimental details to reproduce Figure 6.

1283 **Data generation** For the linearly separable distribution we set $d = 1000$, $n_{\text{train}} = 10^4$, $n_{\text{test}} = 10^5$, $\gamma = 6$.

1285 **Model and hyper-parameters** For all the experiments, we use the one hidden layer architecture defined in Equation (7)
1286 with 100 neurons. We use PyTorch SGD optimiser and train all networks for 100 epochs. We sweep over the learning rate
1287 $\eta \in \{0.1, 0.01, 0.001\}$ and for each perturbation budget, we choose the one that interpolates the training set and minimises
1288 robust error on the test set.

1290 **Robust evaluation** We perform all the attacks to evaluate robust risk at test-time using exact line search; this is computationally tractable since the attacks are directed along one dimension.

1293 **Training paradigms** For standard training (ST), we train the network to minimise the cross-entropy loss. For adversarial
1294 training (AT) (Madry et al., 2018; Goodfellow et al., 2015), we train the network to minimise the robust binary cross-entropy
1295 loss. At each epoch, we compute an exact adversarial example using line search and update the weights using a gradient
1296 with respect to this example. For convex outer adversarial polytope (COAP) (Wong & Kolter, 2018; Wong et al., 2018), at
1297 each epoch, we compute upper and lower bounds u and ℓ as described in Theorem 4. We then train the network to minimise
1298 the upper bound on robust error from Theorem 3.

1300 D.2. Synthetic experiments with ℓ_2 -ball perturbations

1302 Below we provide complete experimental details to reproduce Figures 3 and 4.

1304 **Data generation** For the spheres dataset, we generate a random $x \in \mathbb{R}^d$ where $\|x\|_2$ is either R_1 or R_{-1} , with equal
1305 probability assigned to each norm. We associate with each x a label y such that $y = -1$ if $\|x\|_2 = R_{-1}$ and $y = 1$ if
1306 $\|x\|_2 = R_1$. We can sample uniformly from this distribution by sampling $z \sim \mathcal{N}(0, I_d)$ and then setting $x = \frac{z}{\|z\|_2}R_{-1}$ or
1307 $x = \frac{z}{\|z\|_2}R_1$. For the linearly separable distribution we set $d = 1000$, $n = 50$, $n_{\text{test}} = 10^5$, $\gamma = 6$. For the concentric
1308 spheres distribution we set $d = 100$, $n = 50$, $n_{\text{test}} = 10^5$, $\gamma_{\min} = 1$ and $\gamma_{\max} = 12$.

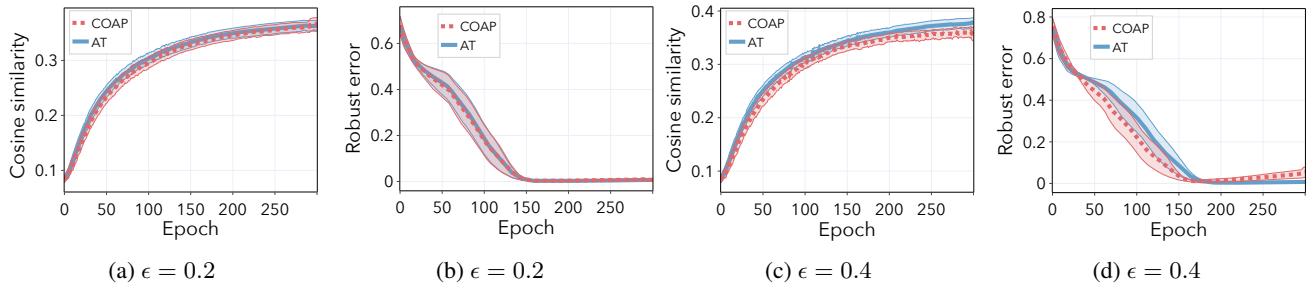
1310 **Model and hyper-parameters** For all the experiments, we use a MLP architecture with $W = 100$ neurons in each hidden
1311 layer and ReLU (\cdot) activation functions. We use PyTorch SGD optimiser with a momentum of 0.95 and train the network
1312 for 150 epochs. We sweep over the learning rate $\eta \in \{0.1, 0.01, 0.001\}$ and for each perturbation budget, we choose the one
1313 that minimises robust error on the test set and interpolates the training set.

1315 **Robust evaluation** We consider ℓ_2 -ball perturbations. We evaluate robust error at test-time using Auto-PGD (Croce &
1316 Hein, 2020) with 100 iterations and 5 random restarts. We use both the cross-entropy and difference of logits loss to prevent
1317 gradient masking. We use the implementation provided in AutoAttack (Croce & Hein, 2020) with minor adjustments to
1318 allow for non-image inputs.

1320 **Training paradigms** For standard training (ST), we train the network to minimise the cross-entropy loss. For adversarial
 1321 training (AT) (Madry et al., 2018; Goodfellow et al., 2015) we train the network to minimise the robust cross-entropy loss.
 1322 At each epoch, we search for adversarial examples using Auto-PGD (Croce & Hein, 2020) with a budget of 10 steps and 1
 1323 random restart. Then, we update the weights using a gradient with respect to the adversarial examples. For convex outer
 1324 adversarial polytope (COAP) (Wong & Kolter, 2018; Wong et al., 2018). We train the network to minimise the upper-bound
 1325 on the robust error. Our implementation is based on the code released by the authors.
 1326

1327 D.3. Additional ℓ_2 -ball perturbation synthetic experiments with small perturbation budget

1328 For completeness we report here the same plots as in Figures 4b and 4c but for smaller perturbation budgets.



1340 *Figure 9.* We report mean and standard error over 15 seeds. Observe that the robust generalisation gap increases with increasing
 1341 perturbation budget ϵ .
 1342

1343 D.4. Image experiments with ℓ_2 -ball perturbations

1344 Below we provide complete experimental details to reproduce Figures 2, 7, 8a and 8b.

1345 **Model architectures** For MNIST, we train the CNN architecture with four convolutional layer and two fully connected
 1346 layers of 512 units introduced in Wong et al. (2018). We report the architectural details in Table 1. For CIFAR-10, we train
 1347 the residual network (ResNet) with the same structure used in Wong et al. (2018); we use 1 residual block with 16, 16, 32,
 1348 and 64 filters. For Tiny ImageNet, we train a WideResNet. Following Xu et al. (2020) we use 3 wide basic blocks with a
 1349 widen factor of 10.

CNN
CONV 32 $3 \times 3 + 1$
CONV 32 $4 \times 4 + 2$
CONV 64 $3 \times 3 + 1$
CONV 64 $4 \times 4 + 2$
FC 512
FC 512

1361 *Table 1.* MNIST model architecture. All layers are followed by ReLU (\cdot) activations. The last fully connected layer is omitted. "CONV k
 1362 $w \times h + s$ " corresponds to a 2D convolutional layer with k filters of size $w \times h$ using a stride of s in both dimensions. "FC n " corresponds
 1363 to a fully connected layer with n outputs.

1364 **Dataset preprocessing** For MNIST, we use full 28×28 images without any augmentations and normalisation. For
 1365 CIFAR-10, we use random horizontal flips and random crops as data augmentation, and normalise images according to
 1366 per-channel statistics. For Tiny ImageNet, we use random crops of 56×56 and random flips during training. During testing,
 1367 we use a central 56×56 crop. We also normalise images according to per-channel statistics.
 1368

1369 **Robust evaluation** We consider ℓ_2 -ball perturbations. We evaluate the robust error using the most expensive version of
 1370 AutoAttack (AA+) (Croce & Hein, 2020). Specifically, we include the following attacks: untargeted APGD-CE (5 restarts),
 1371 untargeted APGD-DLR (5 restarts), untargeted APGD-DLR (5 restarts), Square Attack (5000 queries), targeted APGD-DLR
 1372 (9 target classes) and targeted FAB (9 target classes).

1375 **AT training details** For MNIST, we train 100 epochs using Adam optimiser (Kingma & Ba, 2015) with a learning rate
 1376 of 0.001, momentum of 0.9 and a batch size of 128; we reduce the learning rate by a factor 0.1 at epochs 40 and 80. For
 1377 CIFAR-10 with ResNet, we train 150 epochs using SGD with a learning rate of 0.05 and a batch size of 128; we reduce the
 1378 learning rate by a factor 0.1 at epochs 80 and 120. For Tiny Imagenet and CIFAR-10 with Wide-Resnet we train 200 epochs
 1379 using SGD with a learning rate of 0.1 and a batch size of 512; we reduce the learning rate by a factor 0.1 at epochs 100
 1380 and 150. For the inner optimisation of all models and datasets, adversarial examples are generated with 10 iterations of
 1381 Auto-PGD (Croce & Hein, 2020).

1382
 1383 **COAP training details** We follow the settings proposed by the authors and report them here. For MNIST, we use the
 1384 Adam optimiser (Kingma & Ba, 2015) with a learning rate of 0.001 and a batch size of 50. We schedule ϵ starting from
 1385 0.01 to the desired value over the first 20 epochs, after which we decay the learning rate by a factor of 0.5 every 10 epochs
 1386 for a total of 60 epochs. For CIFAR-10, we use the SGD optimiser with a learning rate of 0.05 and a batch size of 50. We
 1387 schedule ϵ starting from 0.001 to the desired value over the first 20 epochs, after which we decay the learning rate by a factor
 1388 of 0.5 every 10 epochs for a total of 60 epochs. For all datasets and models, we use random projection of 50 dimensions.
 1389 For all experiments, we use the implementation provided in Wong et al. (2018).

1390
 1391 **CROWN-IBP training details** We follow the settings proposed by the authors and report them here. For MNIST, we train
 1392 200 epochs with a batch size of 256. We use Adam optimiser (Kingma & Ba, 2015) and set the learning rate to 5×10^{-4} .
 1393 We warm up with 10 epochs of regular training, and gradually ramp up ϵ_{train} from 0 to ϵ in 50 epochs. We reduce the learning
 1394 rate by a factor 0.1 at epoch 130 and 190. For CIFAR-10, we train 2000 epochs with a batch size of 256, and a learning
 1395 rate of 5×10^{-4} . We warm up for 100 epochs, and ramp-up ϵ for 800 epochs. Learning rate is reduced by a factor 0.1 at
 1396 epoch 1400 and 1700. For Tiny ImageNet, we train 600 epochs with batch size 128. The first 100 epochs are clean training,
 1397 then we gradually increase ϵ_{train} with a schedule length of 400. For all datasets, an hyper-parameter β to balance LiRPA
 1398 bounds and IBP bounds for the output layer is gradually decreased from 1 to 0 (1 for only using LiRPA bounds and 0 for
 1399 only using IBP bounds), with the same schedule of ϵ . For all experiments, we use the implementation provided in the auto
 1400 LiRPA library (Xu et al., 2020).

1402 D.5. Image experiments with signal-directed perturbations

1403 Below we provide complete experimental details to reproduce Figures 5a and 5b.

1404 First, we explain our extensions of COAP to the threat model introduced in Equation (3). Rather than deriving the dual
 1405 problem as in Appendix A.1, we consider the conjugate function view introduced in Wong et al. (2018). In particular, we
 1406 only have to modify the dual of the input layer to the network. Below we derive the conjugate bound for the signal-directed
 1407 threat model:

$$\begin{aligned}
 \sup_{\tilde{x} \in \mathcal{B}(x)} \nu_1^\top \tilde{x} &= \sup_{k, \beta} \nu_1^\top (x + s_k \beta) \\
 &= \nu_1^\top x + \epsilon \max_k |\nu_1^\top s_k|
 \end{aligned}$$

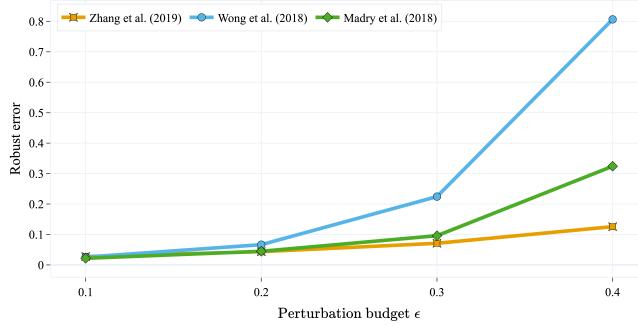
1408 For all experiments, we use the convolutional neural network architecture described in Table 1. Note that it is not possible
 1409 to scale to ResNet with the threat model in Equation (3), as the random projections trick derived in Wong et al. (2018) is
 1410 tailored to ℓ_∞ and ℓ_2 threat models.

1411 **AT training details** For both MNIST and CIFAR-10, we train 20 epochs using Adam optimiser (Kingma & Ba, 2015)
 1412 with a learning rate of 0.001, momentum of 0.9 and a batch size of 64; we reduce the learning rate by a factor 0.1 at epochs
 1413 10. For the inner optimisation of all models and datasets, we solve the exact problem as it is computationally efficient to
 1414 line-search the maximal perturbation.

1415 **COAP training details** For both MNIST and CIFAR-10, we use the Adam optimiser (Kingma & Ba, 2015) with a learning
 1416 rate of 0.001 and a batch size of 64. We schedule ϵ starting from 0.01 to the desired value over the first 3 epochs, after which
 1417 we decay the learning rate by a factor of 0.5 every 10 epochs. For all datasets and models, we do not use random projections.
 1418 For all experiments, we use the implementation provided in Wong et al. (2018).

1430 E. Additional experiments with ℓ_∞ -ball perturbations

1431 We consider ℓ_∞ adversaries and investigate the effect of perturbation budget on robust generalisation. We compare two
 1432 probable defences, CROWN-IBP (Zhang et al., 2020) and COAP (Wong & Kolter, 2018), with the most effective empirical
 1433 defence to date, AT (Madry et al., 2018). We plot the robust error on MNIST for increasing perturbation budget in Figure
 1434 10. As expected, the robustness of COAP quickly degrades as the perturbation budget increases. Especially for larger
 1435 model size and perturbation budget, we observe that AT attains much better robust accuracy than COAP. However, for
 1436 ℓ_∞ adversaries, interval bound propagation methods as CROWN-IBP achieve state of the art performance, outperforming
 1437 empirical defences.



1450 *Figure 10.* Results for ℓ_∞ adversaries against models trained using provable defenses and adversarial training on MNIST. In (a),(b) we
 1451 plot the robust error for a CNN network as the perturbation budget ϵ increases.

1452 In Table 2 we report the robust error on CIFAR10 for CNN and ResNet. For all perturbation budgets considered, AT achieves
 1453 the best robust and standard error. Notably, the standard error of COAP and CROWN-IBP degrades significantly for larger
 1454 perturbation budgets.

METHOD	MODEL	PERTURBATION BUDGET (ϵ)	STANDARD ERROR	ROBUST ERROR
AT (MADRY ET AL., 2018)	CNN	2/255	17.05%	41.52%
AT (MADRY ET AL., 2018)	RESNET	2/255	16.58%	39.52%
COAP (WONG & KOLTER, 2018)	CNN	2/255	32.30%	41.23%
COAP (WONG & KOLTER, 2018)	RESNET	2/255	33.67%	42.75%
CROWN-IBP (ZHANG ET AL., 2020)	CNN	2/255	41.67%	52.71%
AT (MADRY ET AL., 2018)	CNN	8/255	25.01%	74.20%
AT (MADRY ET AL., 2018)	RESNET	8/255	22.59%	69.48%
COAP (WONG & KOLTER, 2018)	CNN	8/255	80.99%	83.54%
COAP (WONG & KOLTER, 2018)	RESNET	8/255	72.93%	77.54%
CROWN-IBP (ZHANG ET AL., 2020)	CNN	8/255	59.91%	71.09%

1469 *Table 2.* Results for ℓ_∞ adversaries against models trained using provable defenses and adversarial training on CIFAR10.