

Adversarial Training for Probabilistic Robustness

Yi Zhang

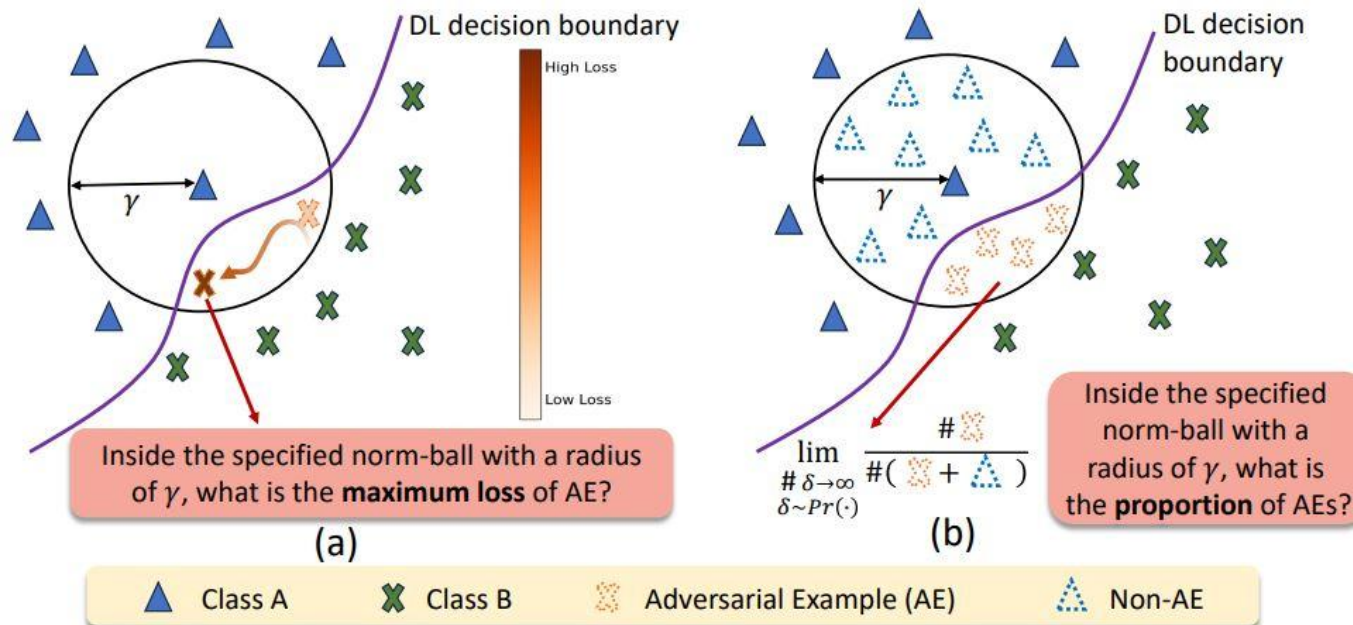
PhD Student within Verification & Validation
WMG, University of Warwick, UK



2025-08-30



Worst-case vs. Probabilistic Robustness



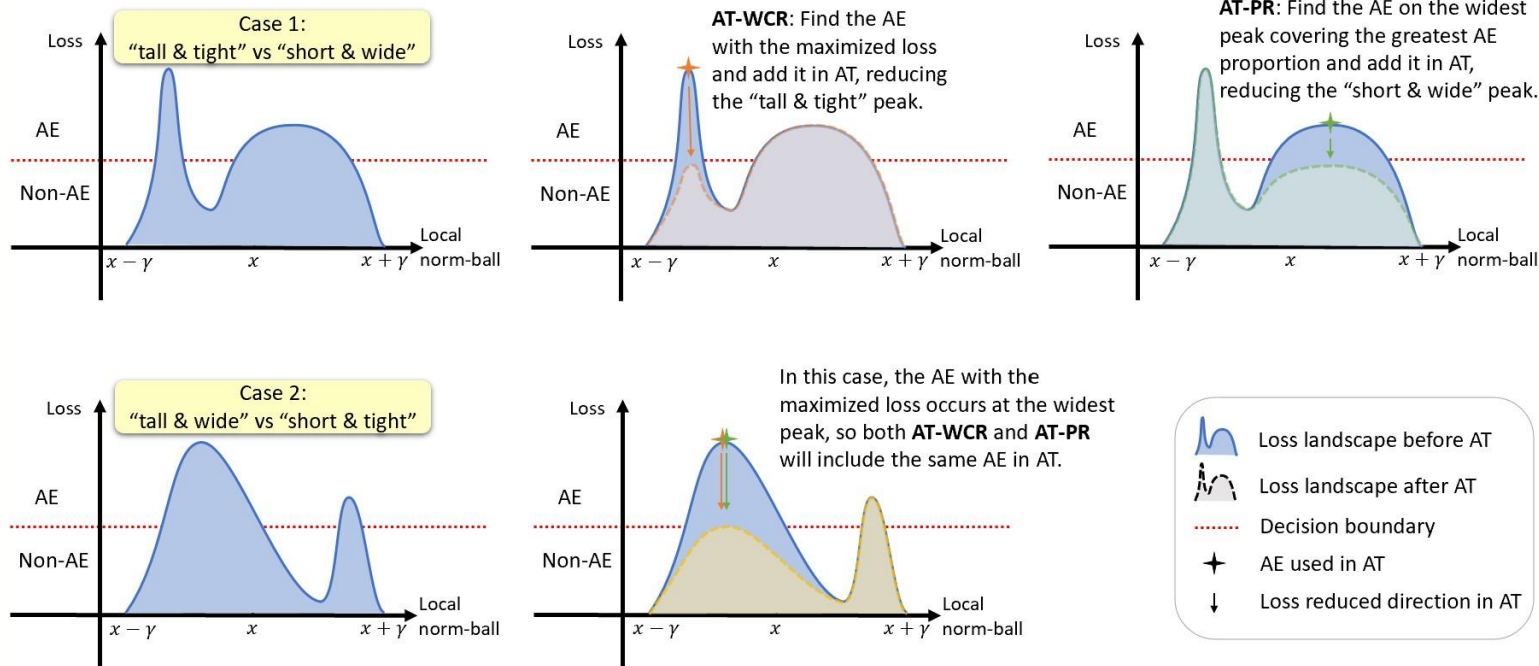
Definition 1 (Probabilistic Robustness) For a DL classifier f_{θ} that takes input x and returns a prediction label, the PR of an input x in a norm ball of radius γ is:

$$PR(x, \gamma) = \mathbb{E}_{\delta \sim Pr(\cdot | x)} [I_{\{f_{\theta}(x + \delta) = y\}}(x + \delta)],$$

$\|\delta\| \leq \gamma$

where $I_S(x)$ is an indicator function—it is equal to 1 when S is true and 0 otherwise; $Pr(\cdot)$ is the local distribution of inputs representing how perturbations δ are generated, which is precisely the “input model” used by [59, 72].

Loss landscape



- Two cases of the local loss landscape and their AEs identified by AT-WCR and AT-PR for AT.

Algorithms of AT-PR

Definition 2 (AT for PR) For a DL model f_{θ} , trained on a dataset \mathcal{D} consisting of pairs (x, y) , the objective of AT-PR can be formulated as:

$$\min_{\theta} \mathbb{E}_{(x,y) \sim \mathcal{D}} \left[\max_{\|\delta\| \leq \gamma, PR(x+\delta, k)=0} k \right],$$

The intuition behind the maximization is to find an optimal δ^* that produces an AE $x' = x + \delta^*$ and to maximize the radius k of a smaller norm-ball centered on x' such that the PR within this k -sized norm-ball is 0 (i.e., all inputs within this smaller norm-ball are AEs).

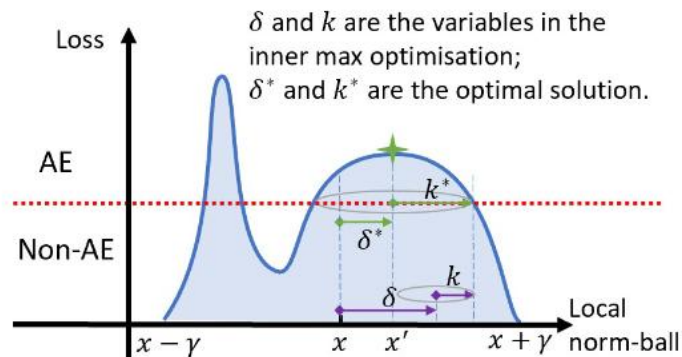


Illustration of the new inner-max optimization in Def 2.

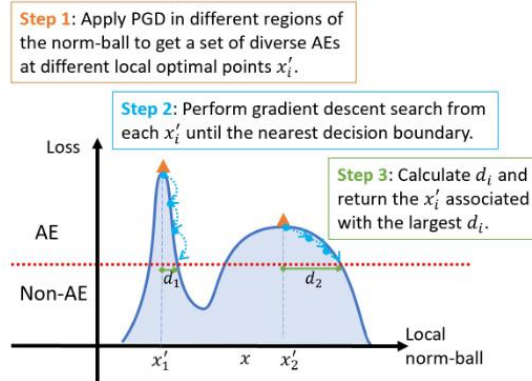


Illustration of the 3 main steps of Algorithm. 1.

Algorithm 1 PGD and gradient-based search for x'_{pr}

Require: Inputs $[X, Y], N, \alpha_{min}, \alpha_{max}, step_{min}, step_{max}$

```

1: Initialize  $AE_s \leftarrow []$   $\triangleright$  Initialize AE candidate set
2: Apply PGD to get different AE candidates
3: for idx = 0 to N do
4:    $x_{init} \leftarrow x + \text{Uniform}(-\gamma, \gamma)$ 
5:    $\alpha \leftarrow \text{Uniform}(\alpha_{min}, \alpha_{max})$ 
6:    $steps \leftarrow \text{Uniform}(step_{min}, step_{max})$ 
7:    $x'_{idx} \leftarrow pgd(x_{init}, y, \gamma, \alpha, steps)$ 
8:   Append  $x'_{idx}$  to  $AE_s$ 
9: end for

```

Require: AE_s, x, y, C

```

10: Initialize  $Max\_d \leftarrow 0$ ;  $x'_{pr} \leftarrow \text{None}$ 
11: for each  $x' \in AE_s$  do
12:    $\tilde{x} = x'$ 
13:   while iter < C do
14:      $y' = f(\tilde{x})$ 
15:     if  $y' = y$  then
16:       break  $\triangleright$  Exit if  $\tilde{x}$  is classified correctly
17:     end if
18:      $g \leftarrow \nabla_{\tilde{x}} L(\tilde{x}, y)$   $\triangleright$  Compute gradient
19:      $\tilde{x} = \tilde{x} - \alpha \cdot g$   $\triangleright$  Update  $\tilde{x}$ 
20:   end while
21:    $d \leftarrow D(\tilde{x}, x')$   $\triangleright$  Distance between  $\tilde{x}$  &  $x'$ 
22:   if  $d > Max\_d$  then
23:      $x'_{pr} = x'$ 
24:      $Max\_d = d$ 
25:   end if
26: end for
27: Output  $x'_{pr}$   $\triangleright$  Farthest AE from decision boundary.

```


Experiment of AT-PR

Table 1. Evaluation results of DL models trained with different AT schemes, showing standard accuracy, adversarial accuracy (indicating WCR: PGD-10, C&W-10), and PR estimated (mean \pm std) for four norm-ball sizes γ . All results are calculated over the test datasets.

Dataset	Model	AT Scheme	Accu. %	PGD Ac. %	CW Ac. %	PR(γ , 0.08) %	PR(γ , 0.1) %	PR(γ , 0.12) %	PR(γ , 0.15) %
CIFAR-10	ResNet-18	W/O AT	94.96	0.02	0.51	78.38 \pm 33.14	66.3 \pm 38.06	53.41 \pm 40.76	37.01 \pm 39.77
		AT-WCR(FGSM)	88.6	0.09	0.09	93.99 \pm 17.5	88.02 \pm 25.03	80.42 \pm 31.81	67.77 \pm 38.6
		AT-WCR(PGD)	81	45.01	44.56	97.28 \pm 13.29	95.77 \pm 17.07	94.1 \pm 20.2	89.78 \pm 26.53
		AT-PR	80.14	39.77	39.03	98.22 \pm 10.38	97.28 \pm 13.16	96.05 \pm 16.18	93.21 \pm 21.47
	WRN-50-2	W/O AT	94.14	0.23	0.14	76.96 \pm 32.76	65.65 \pm 37.41	54.43 \pm 39.88	40.19 \pm 39.49
		AT-WCR(FGSM)	92.04	0.08	0.04	95.68 \pm 14.53	91.9 \pm 20.25	86.67 \pm 25.98	76.79 \pm 33.1
CIFAR-100	ResNet-18	AT-WCR(PGD)	82.62	49.51	48.13	97.86 \pm 12.61	97.08 \pm 14.54	95.84 \pm 17.14	92.62 \pm 22.87
		AT-PR	82.71	43.27	43.07	97.33 \pm 13.03	96.22 \pm 15.58	93.56 \pm 20.89	
	WRN-50-2	W/O AT	75.5	0.02	0.0	44.58 \pm 42.64	32.1 \pm 40.1	23.32 \pm 36.41	15.62 \pm 31.57
		AT-WCR(FGSM)	67.80	1.41	0.66	83.90 \pm 31.11	72.89 \pm 37.94	61.06 \pm 42.23	44.71 \pm 43.04
	ResNet-18	AT-WCR(PGD)	52.85	19.93	19.51	95.4 \pm 18.15	92.7 \pm 22.72	88.98 \pm 28.19	80.47 \pm 36.4
		AT-PR	53.21	19.38	18.39	97.59 \pm 12.96	95.96 \pm 16.62	93.15 \pm 21.77	87.38 \pm 29.63
SVHN	WRN-50-2	W/O AT	75.46	0.49	0.45	51.11 \pm 40.54	37.96 \pm 39.23	28.33 \pm 36.72	19.11 \pm 31.97
		AT-WCR(FGSM)	73.17	0.29	0.1	85.68 \pm 27.45	76.62 \pm 33.89	66.37 \pm 38.66	52.04 \pm 41.42
		AT-WCR(PGD)	57.46	23.17	22.5	95.31 \pm 17.38	92.0 \pm 23.79	88.14 \pm 29.15	80.67 \pm 36.21
		AT-PR	59.46	22.06	21.59	96.05 \pm 16.17	93.72 \pm 20.47	89.48 \pm 26.87	81.6 \pm 34.86
	ResNet-18	W/O AT	94.96	0.08	0.12	97.11 \pm 9.66	95.63 \pm 11.98	93.74 \pm 14.48	90.12 \pm 18.40
		AT-WCR(FGSM)	93.49	4.68	5.32	98.79 \pm 7.06	98.42 \pm 7.70	97.86 \pm 8.67	96.57 \pm 10.78
Tiny-ImageNet	WRN-50-2	AT-WCR(PGD)	89.21	44.57	43.70	98.96 \pm 5.85	98.50 \pm 7.40	97.72 \pm 9.61	95.41 \pm 14.41
		AT-PR	90.29	45.81	42.21	99.12 \pm 5.15	98.87 \pm 5.93	98.54 \pm 6.84	97.84 \pm 8.60
	ResNet-18	W/O AT	95.91	1.74	1.29	96.60 \pm 10.36	94.77 \pm 13.01	92.40 \pm 15.87	87.97 \pm 20.04
		AT-WCR(FGSM)	94.31	8.13	8.54	98.34 \pm 7.49	97.71 \pm 8.39	96.80 \pm 9.69	95.12 \pm 12.06
	WRN-50-2	AT-WCR(PGD)	88.91	54.99	50.58	98.72 \pm 6.44	98.28 \pm 7.66	97.76 \pm 8.95	96.77 \pm 11.18
		AT-PR	91.97	52.22	51.58	98.88 \pm 6.12	98.38 \pm 7.53	97.72 \pm 9.19	96.26 \pm 12.35
SVHN	WRN-50-2	W/O AT	59.26	0.0	0.0	80.78 \pm 33.1	72.95 \pm 37.69	64.17 \pm 40.91	51.14 \pm 42.91
		AT-WCR(FGSM)	54.81	0.05	0.01	89.9 \pm 24.97	85.33 \pm 29.95	80.18 \pm 33.74	69.03 \pm 39.39
		AT-WCR(PGD)	47.43	21.3	19.31	96.32 \pm 17.01	93.79 \pm 22.1	90.16 \pm 27.76	82.46 \pm 35.97
		AT-PR	51.61	21.1	19.89	98.23 \pm 11.06	96.1 \pm 17.84	93.12 \pm 23.32	86.65 \pm 32.26
	ResNet-18	W/O AT	61.66	0.0	0.0	82.39 \pm 31.3	75.2 \pm 36.0	66.95 \pm 39.42	53.7 \pm 42.17
		AT-WCR(FGSM)	57	0.13	0.02	90.67 \pm 23.92	85.81 \pm 29.5	80.57 \pm 34.13	70.22 \pm 39.36
Tiny-ImageNet	WRN-50-2	AT-WCR(PGD)	49.33	20.97	19.58	95.45 \pm 18.55	93.08 \pm 23.48	89.85 \pm 28.3	83.53 \pm 34.99
		AT-PR	53.98	21.72	20.95	96.71 \pm 15.74	95.02 \pm 19.89	92.5 \pm 23.87	86.23 \pm 32.29
	ResNet-34	W/O AT	61.66	0.0	0.0	82.39 \pm 31.3	75.2 \pm 36.0	66.95 \pm 39.42	53.7 \pm 42.17
		AT-WCR(FGSM)	57	0.13	0.02	90.67 \pm 23.92	85.81 \pm 29.5	80.57 \pm 34.13	70.22 \pm 39.36
	WRN-50-2	AT-WCR(PGD)	49.33	20.97	19.58	95.45 \pm 18.55	93.08 \pm 23.48	89.85 \pm 28.3	83.53 \pm 34.99
		AT-PR	53.98	21.72	20.95	96.71 \pm 15.74	95.02 \pm 19.89	92.5 \pm 23.87	86.23 \pm 32.29

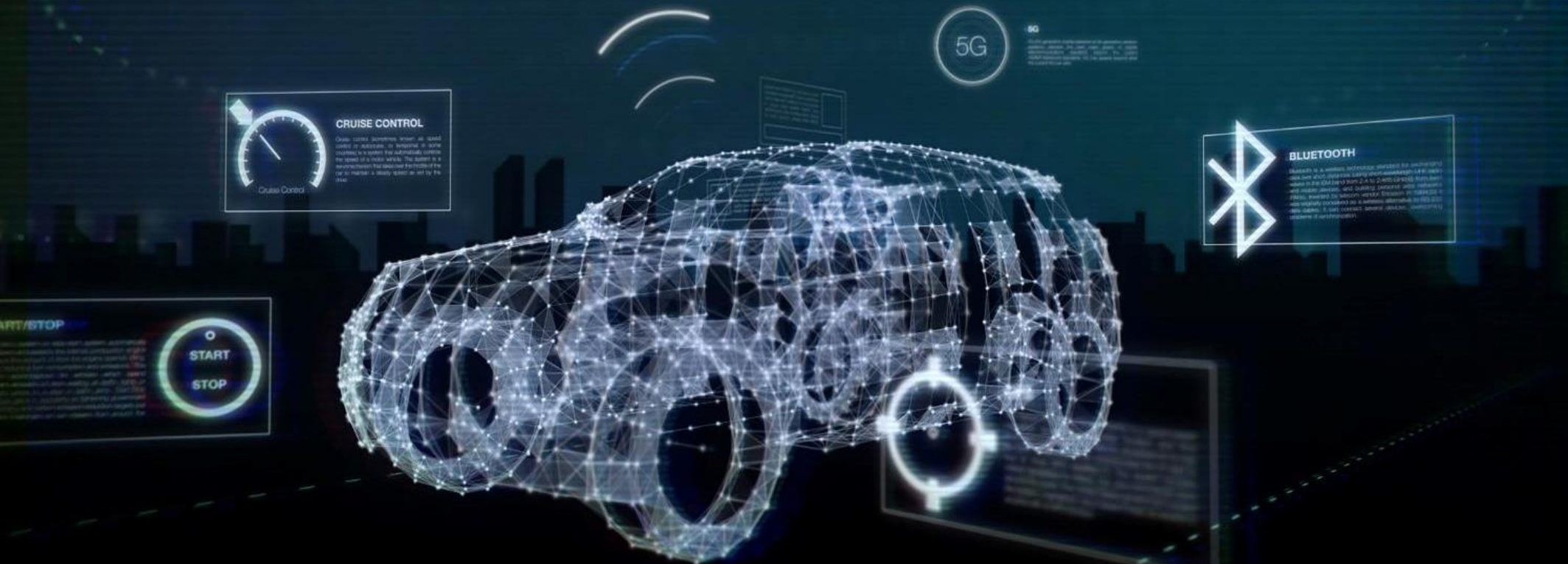
Table 2. ViT models trained on CIFAR-10 with TRADES and AT-PR, showing standard accuracy, adversarial accuracy (indicating WCR: PGD, C&W, AA), and PR estimated (mean \pm std) for four norm-ball sizes γ .

Dataset	Model	AT Scheme	Accu. %	PGD-20. %	CW-20. %	AA. %	PR(γ , 0.08) %	PR(γ , 0.1) %	PR(γ , 0.12) %	PR(γ , 0.15) %
CIFAR-10	DeiT-Ti	TRADES	79.44	49.78	46.45	45.96	96.61 \pm 15.23	94.73 \pm 19.36	92.69 \pm 23.04	89.18 \pm 27.9
		+AT-PR	81.07	49.08	46.92	45.95	97.26 \pm 14.0	95.94 \pm 17.12	94.16 \pm 20.5	90.65 \pm 25.88
	ViT-S	TRADES	81.03	50.96	48.14	47.40	96.93 \pm 14.99	95.48 \pm 18.4	93.58 \pm 21.62	90.06 \pm 26.81
		+AT-PR	83.29	50.44	48.20	47.34	97.76 \pm 11.89	96.03 \pm 16.46	93.98 \pm 20.69	90.48 \pm 26.19
	ViT-B	TRADES	83.55	53.14	50.94	49.96	97.34 \pm 13.57	95.87 \pm 16.99	93.78 \pm 21.17	90.77 \pm 25.88
		+AT-PR	86.03	52.94	51.24	49.34	97.9 \pm 11.8	96.49 \pm 15.2	94.44 \pm 19.62	91.02 \pm 25.27

Table 3. Comparison of PR(γ , γ) and training time (sec/epoch) for more recent AT methods (TRADES, MART, ALP, CLP) and Fast-AT-PR

Data	Model	Metric	AT Scheme						
			PGD	TRADES	MART	ALP	CLP	Fast-AT-PR	AT-PR
CIFAR-10	ResNet-18	PR(γ , 0.08) %	97.28 \pm 13.29	97.74 \pm 12.42	97.4 \pm 13.7	96.73 \pm 15.9	97.61 \pm 13.06	98.12 \pm 11.52	98.22 \pm 10.38
		PR(γ , 0.10) %	95.77 \pm 17.07	96.33 \pm 16.39	95.92 \pm 17.43	95.16 \pm 19.53	96.29 \pm 16.66	97.1 \pm 14.34	97.28 \pm 13.16
		PR(γ , 0.12) %	94.10 \pm 20.20	94.55 \pm 20.31	93.93 \pm 21.38	93.24 \pm 23.08	94.54 \pm 20.28	95.4 \pm 17.90	96.05 \pm 16.18
		PR(γ , 0.15) %	89.78 \pm 26.53	91.22 \pm 25.70	90.35 \pm 26.91	89.76 \pm 28.14	90.86 \pm 26.23	92.71 \pm 23.54	93.21 \pm 21.47
	Time (s/ep.)	54	61	44	64	64	97	184	
	WRN-50-2	PR(γ , 0.08) %	97.86 \pm 12.61	97.55 \pm 13.14	97.40 \pm 13.61	97.75 \pm 12.43	97.40 \pm 13.91	97.96 \pm 10.82	98.09 \pm 10.81
PR(γ , 0.10) %		97.08 \pm 14.54	96.14 \pm 16.70	96.15 \pm 16.78	96.42 \pm 16.11	96.20 \pm 17.58	97.22 \pm 12.96	97.33 \pm 13.03	
PR(γ , 0.12) %		95.84 \pm 17.14	94.33 \pm 20.41	94.41 \pm 20.37	94.64 \pm 20.01	94.95 \pm 19.88	96.22 \pm 15.62	96.22 \pm 15.58	
PR(γ , 0.15) %		92.62 \pm 22.87	91.01 \pm 25.79	91.11 \pm 25.65	91.22 \pm 25.70	91.48 \pm 25.69	94.74 \pm 19.88	93.56 \pm 20.89	
CIFAR-100	ResNet-18	Time (s/ep.)	131	105	71	155	157	229	305
		PR(γ , 0.08) %	95.40 \pm 18.15	95.02 \pm 19.27	95.45 \pm 18.11	96.20 \pm 16.07	96.16 \pm 16.36	97.09 \pm 13.48	97.59 \pm 12.96
		PR(γ , 0.10) %	92.70 \pm 22.72	92.09 \pm 24.22	92.73 \pm 23.39	94.40 \pm 20.08	93.81 \pm 21.42	94.67 \pm 18.91	95.96 \pm 16.62
		PR(γ , 0.12) %	88.98 \pm 28.19	87.80 \pm 29.80	89.35 \pm 28.32	92.05 \pm 24.12	91.05 \pm 25.86	92.63 \pm 23.91	93.15 \pm 21.77
	PR(γ , 0.15) %	80.47 \pm 36.40	79.20 \pm 37.73	81.93 \pm 35.72	85.89 \pm 31.65	84.97 \pm 33.00	85.87 \pm 31.99	87.38 \pm 29.63	
	Time (s/ep.)	55	61	43	64	65	97	180	
WRN-50-2	PR(γ , 0.08) %	95.31 \pm 17.38	95.24 \pm 18.61	95.30 \pm 16.87	95.34 \pm 17.79	95.29 \pm 17.70	96.52 \pm 15.11	96.05 \pm 16.17	
	PR(γ , 0.10) %	92.00 \pm 23.79	92.30 \pm 23.61	92.59 \pm 22.81	92.36 \pm 23.19	91.79 \pm 23.87	93.49 \pm 20.93	93.72 \pm 20.47	
	PR(γ , 0.12) %	88.14 \pm 29.15	88.65 \pm 28.46	88.38 \pm 28.65	88.51 \pm 28.22	87.54 \pm 29.92	88.82 \pm 28.38	89.48 \pm 26.87	
	PR(γ , 0.15) %	80.67 \pm 36.21	80.35 \pm 36.54	80.79 \pm 36.01	80.37 \pm 36.86	79.71 \pm 36.76	81.02 \pm 37.09	81.60 \pm 34.86	
SVHN	ResNet-18	Time (s/ep.)	131	102	70	155	156	225	312
		PR(γ , 0.08) %	98.96 \pm 5.85	99.03 \pm 5.76	98.99 \pm 6.15	99.05 \pm 5.73	99.12 \pm 5.33	98.90 \pm 5.97	99.12 \pm 5.15
		PR(γ , 0.10) %	98.50 \pm 7.40	98.46 \pm 8.27	98.31 \pm 8.20	98.63 \pm 7.15	98.60 \pm 7.21	98.53 \pm 7.24	98.87 \pm 5.93
		PR(γ , 0.12) %	97.72 \pm 9.61	97.30 \pm 11.08	97.10 \pm 11.26	97.85 \pm 9.13	97.66 \pm 9.78	98.11 \pm 8.19	98.54 \pm 6.84
	PR(γ , 0.15) %	95.41 \pm 14.41	94.15 \pm 16.57	94.12 \pm 16.65	95.54 \pm 14.10	94.93 \pm 15.00	97.10 \pm 10.05	97.84 \pm 8.60	
	Time (s/ep.)	78	86	61	91	93	141	202	
WRN-50-2	PR(γ , 0.08) %	98.72 \pm 6.44	98.69 \pm 6.90	98.61 \pm 7.66	98.54 \pm 7.52	98.59 \pm 7.46	98.43 \pm 6.69	98.88 \pm 6.12	
	PR(γ , 0.10) %	98.28 \pm 7.66	98.11 \pm 8.66	98.07 \pm 9.27	98.12 \pm 8.74	98.00 \pm 8.89	98.15 \pm 7.29	98.38 \pm 7.53	
	PR(γ , 0.12) %	97.76 \pm 8.95	97.26 \pm 10.89	97.34 \pm 10.90	97.65 \pm 8.85	96.98 \pm 11.13	97.79 \pm 8.12	97.72 \pm 9.19	
	PR(γ , 0.15) %	96.77 \pm 11.18	95.20 \pm 14.84	95.45 \pm 14.72	95.67 \pm 14.32	94.74 \pm 15.37	96.94 \pm 9.80	96.26 \pm 12.35	
Time (s/ep.)	190	146	102	223	225	321	420		

Thank you...



Yi Zhang

Yi.Zhang.16@warwick.ac.uk