

A. Dense Adversary Generation attack

In this section, we provide the algorithm of the Dense Adversary Generation (DAG) attack from [39]. The main difference with the original algorithm proposed in [39] is the stopping criterion based on the pixel success rate in steps 9 to 11. In the original method published in [39], the attack is supposed to stop once all the constraints are satisfied (*i.e.* all pixels in the mask \mathbf{m} are adversarial). However, since this criterion is rarely satisfied, even after hundreds of iterations on a dataset like Cityscapes, the actual implementation made available in <https://github.com/cihangxie/DAG> uses the stopping criterion described in **Algorithm 2**: the attack stops once a threshold of pixel success rate set to 99% is reached. The threshold value used is thus identical to the one used in the experiments of this paper.

Algorithm 2 DAG attack

Require: Classifier f , original image $\mathbf{x} \in [0, 1]^{C \times H \times W}$, true or target label $\mathbf{y} \in \mathbb{N}^{H \times W}$, binary mask $\mathbf{m} \in \{0, 1\}^{H \times W}$

Require: Step size η , maximum number of iterations N , threshold of pixel success rate ν

- 1: Initialize $\delta^{(0)} \leftarrow \mathbf{0}$
- 2: If targeted attack: $\mu \leftarrow -1$ else $\mu \leftarrow 1$
- 3: **for** $t \leftarrow 1, \dots, N$ **do**
- 4: $\tilde{\mathbf{x}}^{(t)} \leftarrow \mathcal{P}_{[0,1]}(\mathbf{x} + \delta^{(t-1)})$ // $\in [0, 1]^{C \times H \times W}$
- 5: $\mathbf{z} \leftarrow f(\tilde{\mathbf{x}}^{(t)})$ // $\in \mathbb{R}^{K \times H \times W}$
- 6: **for** $i \leftarrow 1, \dots, d$ **do**
- 7: $\Delta \mathbf{z}_i = \mu(\mathbf{z}_{\mathbf{y}_i, i} - \max_{j \neq \mathbf{y}_i} \mathbf{z}_{j, i})$ // Difference of logits
- 8: **end for**
- 9: $r \leftarrow \frac{\mathbf{m}^\top [\Delta \mathbf{z} < 0]}{\|\mathbf{m}\|_1}$ // Pixel success rate
- 10: **if** $r \geq \nu$ **then**
- 11: **return** $\tilde{\mathbf{x}}^{(t)}$ // Stop the attack
- 12: **end if**
- 13: $\mathcal{L} \leftarrow \mathbf{m}^\top \max\{0, \Delta \mathbf{z}_i\}$
- 14: $\mathbf{g} \leftarrow \nabla_{\delta} \mathcal{L}$
- 15: $\delta^{(t)} \leftarrow \delta^{(t-1)} - \frac{\eta}{\|\mathbf{g}\|_\infty} \mathbf{g}$ // Normalized gradient step
- 16: **end for**

B. Proof of Proposition 1

Proof. Problem (12) amounts to minimizing function

$$\Phi: (\mathbf{p}, \beta) \mapsto \frac{1}{2} \|\mathbf{p} - \delta\|_2^2 + \lambda\beta + \iota_{[0, +\infty[}(\beta \mathbf{1}_{Cd} - \mathbf{p}) + \iota_{[0, +\infty[}(\mathbf{p} + \beta \mathbf{1}_{Cd}) + \iota_\Lambda(\mathbf{p}), \quad (20)$$

defined on $\mathbb{R}^{Cd} \times \mathbb{R}$. This function is convex since it is a sum of elementary convex functions of (\mathbf{p}, β) . It follows that its marginal function

$$\underline{\Phi}: \beta \mapsto \inf_{\mathbf{p} \in \mathbb{R}^{Cd}} \Phi(\mathbf{p}, \beta) \quad (21)$$

is convex. Since, for any given $\beta \in [0, 1]$, $\mathbf{p} \mapsto \Phi(\mathbf{p}, \beta)$ is strongly convex and proper, it admits a unique minimizer \mathbf{p}_β . More precisely, we have, for every $\beta \in \mathbb{R}$,

$$\underline{\Phi}(\beta) = \begin{cases} \frac{1}{2} \|\mathbf{p}_\beta - \delta\|_2^2 + \lambda\beta & \text{if } \beta \in [0, 1] \\ +\infty & \text{otherwise.} \end{cases} \quad (22)$$

The above function admits a unique minimizer $\beta^* \in [0, 1]$ since $\mathbf{p}^* = \mathbf{p}_{\beta^*}$ is uniquely defined (as it is the proximity point of a proper lower-semicontinuous convex function) and $\beta^* = \|\mathbf{p}_{\beta^*}\|_\infty$.

Let $\mathbf{p}^* = \text{prox}_{\lambda\|\cdot\|_\infty + \iota_\Lambda}(\delta)$ be the minimizer of (12) and let $\delta_\Lambda = \mathcal{P}_\Lambda(\delta)$. We have

$$\frac{1}{2} \|\mathbf{p}^* - \delta\|_2^2 + \lambda \|\mathbf{p}^*\|_\infty \leq \frac{1}{2} \|\delta_\Lambda - \delta\|_2^2 + \lambda \|\delta_\Lambda\|_\infty \quad (23)$$

Since $\delta_\Lambda = \mathcal{P}_\Lambda(\delta) = \arg \min_{\mathbf{y} \in \Lambda} \|\mathbf{y} - \delta\|_2^2$, we also have

$$\|\delta_\Lambda - \delta\|_2^2 \leq \|\mathbf{p}^* - \delta\|_2^2 \quad (24)$$

Thus

$$\begin{aligned} \lambda \|\mathbf{p}^*\|_\infty &\leq \underbrace{\frac{1}{2} (\|\delta_\Lambda - \delta\|_2^2 - \|\mathbf{p}^* - \delta\|_2^2)}_{\leq 0} + \lambda \|\delta_\Lambda\|_\infty \\ &\leq \lambda \|\delta_\Lambda\|_\infty \end{aligned} \quad (25)$$

We deduce that $\boxed{\beta^* = \|\mathbf{p}^*\|_\infty \leq \|\delta_\Lambda\|_\infty}$. \square

C. ALMA prox attack algorithm

Algorithm 3 ALMA prox attack (untargeted)

Require: Classifier f , original image $\mathbf{x} \in [0, 1]^{C \times H \times W}$, true or target label $\mathbf{y} \in \mathbb{N}^{H \times W}$, binary mask $\mathbf{m} \in \{0, 1\}^{H \times W}$

Require: Threshold of pixel success rate ν

Require: Penalty function P , initial multiplier $\boldsymbol{\mu}^{(0)} \in \mathbb{R}_{++}^{H \times W}$, initial penalty parameter $\boldsymbol{\rho}^{(0)} \in \mathbb{R}_{++}^{H \times W}$

Require: Minimum scale w_{\min} , scale adjustment rate $\gamma_w > 0$

Require: Number of iterations N , initial step size $\lambda^{(0)}$, penalty parameter increase rate $\gamma > 1$, constraint improvement rate τ , M number of steps between ρ increase, α smoothing parameter

- 1: Initialize $\boldsymbol{\delta}^{(0)} \leftarrow \mathbf{0}$, $\mathbf{v}^{(0)} \leftarrow \mathbf{0}$, $w^{(0)} \leftarrow 1$
- 2: **for** $t \leftarrow 1, \dots, N$ **do**
- 3: $\tilde{\mathbf{x}}^{(t)} \leftarrow \mathbf{x} + \boldsymbol{\delta}^{(t-1)}$ // $\in [0, 1]^{C \times H \times W}$
- 4: $\mathbf{d}^{(t)} \leftarrow \text{DLR}^+(f(\tilde{\mathbf{x}}^{(t)}), \mathbf{y})$ // $\in \mathbb{R}^{H \times W}$
- 5: **if** $\frac{\mathbf{m}^\top [\mathbf{d}^{(t)} \leq 0]}{\|\mathbf{m}\|_1} < \tau$ **then** // Adjust constraint scale
- 6: $\hat{w} \leftarrow \frac{w^{(t-1)}}{1-\gamma_w}$ // Increase scale
- 7: **else**
- 8: $\hat{w} \leftarrow \frac{w^{(t-1)}}{1+\gamma_w}$ // Decrease scale
- 9: **end if**
- 10: $w^{(t)} \leftarrow \mathcal{P}_{[w_{\min}, 1]}(\hat{w})$
- 11: $\xi^{(t)} \leftarrow (1 - (1 - \nu)^{\frac{t-1}{N-1}})$ -percentile of $\mathbf{d}^{(t)}$
- 12: $\tilde{\mathbf{m}}^{(t)} \leftarrow [\mathbf{d}^{(t)} \leq \xi^{(t)}]$ // $\in \{0, 1\}^{H \times W}$
- 13: $\hat{\boldsymbol{\mu}} \leftarrow \nabla_{\mathbf{d}} \left((\tilde{\mathbf{m}}^{(t)})^\top P(w^{(t)} \mathbf{d}^{(t)}, \boldsymbol{\rho}^{(t-1)}, \boldsymbol{\mu}^{(t-1)}) \right)$
- 14: $\boldsymbol{\mu}^{(t)} \leftarrow \mathcal{P}_{[\mu_{\min}, \mu_{\max}]}(\alpha \boldsymbol{\mu}^{(t-1)} + (1 - \alpha) \hat{\boldsymbol{\mu}})$ // $\in \mathbb{R}_{++}^{H \times W}$
- 15: **for** $i \leftarrow 1, \dots, d$ **do** // ρ adjustment
- 16: **if** $t \bmod M = 0$ **and** $\tilde{m}_i^{(t)} = 1$ **and** $(\exists j \in \{0, \dots, M-1\} : \mathbf{d}_i^{(t-j)} \leq 0 \text{ or } \mathbf{d}_i^{(t)} \leq \tau \mathbf{d}_i^{(t-M)})$ **then**
- 17: $\rho_i^{(t)} \leftarrow \rho_i^{(t-1)}$ // Constraint improved or satisfied
- 18: **else**
- 19: $\rho_i^{(t)} \leftarrow \gamma \rho_i^{(t-1)}$
- 20: **end if**
- 21: **end for**
- 22: $\mathcal{L} \leftarrow (\tilde{\mathbf{m}}^{(t)})^\top P(w^{(t)} \mathbf{d}^{(t)}, \boldsymbol{\rho}^{(t)}, \boldsymbol{\mu}^{(t)})$ // $\in \mathbb{R}$
- 23: $\mathbf{g}^{(t)} \leftarrow \nabla_{\boldsymbol{\delta}} \mathcal{L}$ // $\in \mathbb{R}^{C \times H \times W}$
- 24: $\mathbf{v}^{(t)} \leftarrow \alpha \mathbf{v}^{(t-1)} + (1 - \alpha) (\mathbf{g}^{(t)})^2$ // $\in \mathbb{R}_+^{C \times H \times W}$
- 25: $\mathbf{H} \leftarrow \text{Diag} \left(\sqrt{\frac{\mathbf{v}^{(t)}}{1 - \alpha^t}} + \varepsilon \right)$
- 26: $\boldsymbol{\delta}^{(t)} \leftarrow \text{prox}_{\lambda^{(t)} \|\cdot\|_\infty + \iota_\Lambda}^{\mathbf{H}} (\boldsymbol{\delta}^{(t-1)} - \lambda^{(t)} \mathbf{H}^{-1} \mathbf{g})$ // VMFB
- 27: **end for**
- 28: **return** $\tilde{\mathbf{x}}^{(t)}$ that is adversarial and has the smallest norm

D. Image size and performance of the models

For all models, except DeepLabV3 DDC-AT, the images of Pascal VOC 2012 are resized so that the smaller side is of length 512 while keeping the aspect ratio, and for Cityscapes, the images keep their original size of 2048×1024 . For DeepLabV3 DDC-AT from [41], the images of Pascal VOC 2012 are resized so that the longer side is of length 512 while keeping the aspect ratio, and for Cityscapes, the images are resized to 1024×512 .

Dataset	Model	mIoU (%)	Pixel Accuracy (%)
Pascal VOC 2012 (+Aug)	DeepLabV3+ ResNet-50 [11]	77.4 \pm 0.8	94.9 \pm 0.2
	DeepLabV3+ ResNet-101 [11]	78.8 \pm 0.1	95.3 \pm 0.1
	FCN HRNetV2 W48 [38]	76.4 \pm 0.2	94.7 \pm 0.1
	DeepLabV3 DDC-AT [41]	75.2 \pm 0.0	94.4
Cityscapes	DeepLabV3+ ResNet-50 [11]	80.1 \pm 0.2	96.4 \pm 0.1
	FCN HRNetV2 W48 [38]	80.5 \pm 0.2	96.6 \pm 0.1
	SegFormer MiT-B0 [40]	76.4 \pm 0.1	95.9 \pm 0.0
	SegFormer MiT-B3 [40]	81.8 \pm 0.0	96.7 \pm 0.1
	DeepLabV3 DDC-AT [41]	71.0 \pm 0.3	95.0

Table 3. Performance of the models used in the experiments on the validation sets. Numbers were obtained from our evaluation; subscripts correspond to the difference with the original evaluation protocol. For DeepLabV3 DDC-AT from [41], the pixel accuracy was not reported.

E. Cityscapes target label

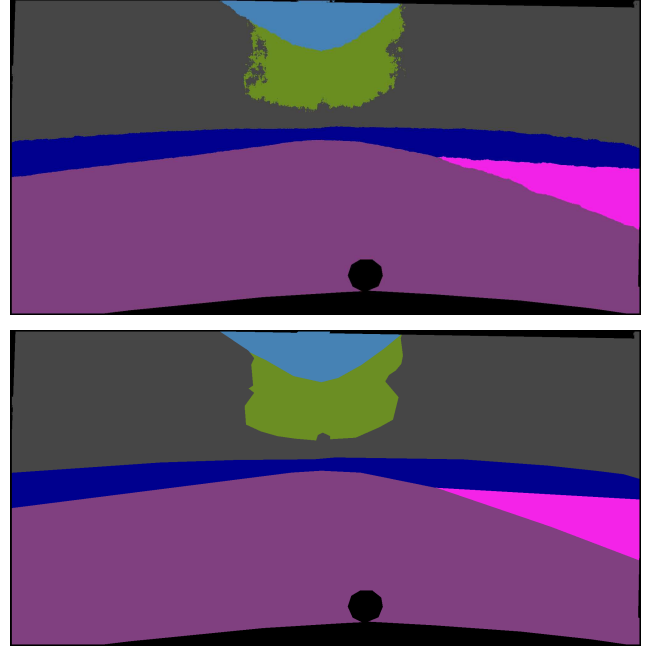


Figure 4. Cityscapes target segmentation used in our experiments. Top is the pixel majority label, bottom is the smoothed version. Classes in the target are: road (purple), sidewalk (pink), car (blue), building (dark gray), vegetation (green), sky (cyan) and no label (black).

F. Masking strategies for ALMA prox and PDPGD

In this section, we detail the modifications brought to PDPGD [30] and study the effect of the masking strategies on PDPGD and ALMA prox. Besides the addition of perturbation tracking logic (*i.e.* monitoring the best smallest perturbations during the optimization), we also incorporate the masking needed for the unlabeled regions.

Modifications for PDPGD For PDPGD, the dual variable $\lambda \in \mathbb{R}$ is replaced by a vector $\boldsymbol{\lambda} \in \mathbb{R}^d$. In [30], the gradient ascent on the dual variable is performed in the log domain, and is projected onto the 1-simplex (section IV of [30]). For the segmentation variant, this translates into adding 1 to the denominator of the softmax function used to project onto the $(d-1)$ -simplex. This is equivalent to padding $\boldsymbol{\lambda}$ with a 0 and projecting it on the d -simplex $\Delta^d \subset [0, 1]^{d+1}$. The weight of the norm $\|\delta\|$ in Equation (10) of [30] becomes 1 minus the sum of the weights of the constraints. Adding the mask \mathbf{m} , this results in the following computations (introducing variables not described in [30]):

$$\begin{aligned} \boldsymbol{\lambda}_\Delta &= \frac{\mathbf{m} \odot \exp \boldsymbol{\lambda}}{1 + \mathbf{m}^\top \exp \boldsymbol{\lambda}} \in [0, 1]^d \\ \mathbb{L}_\delta(\delta, \boldsymbol{\lambda}) &= (1 - \mathbf{m}^\top \boldsymbol{\lambda}_\Delta) \|\delta\| + \boldsymbol{\lambda}_\Delta^\top \mathcal{L}(\mathbf{x} + \delta, \mathbf{y}). \end{aligned} \quad (26)$$

By replacing \mathbf{m} by $\tilde{\mathbf{m}}^{(t)}$ in each iteration, we obtain the adaptive constraint masking described in Section 4.1.

The step-size is set to 0.01 for the primal variables and 0.1 for the dual variables with the same exponential and linear decays respectively, as in [30]. The dual variable is initialized such that the ratio of the weight of the norm term and the constraints terms is 1. From the above equations, we can derive the initial $\boldsymbol{\lambda}^{(0)} \in \mathbb{R}^d$ from the ratio $r \in \mathbb{R}_{++}$:

$$\begin{aligned} \frac{1 - \mathbf{m}^\top \boldsymbol{\lambda}_\Delta^{(0)}}{\mathbf{m}^\top \boldsymbol{\lambda}_\Delta^{(0)}} &= r \Leftrightarrow 1 - \mathbf{m}^\top \boldsymbol{\lambda}_\Delta^{(0)} = r \mathbf{m}^\top \boldsymbol{\lambda}_\Delta^{(0)} \\ &\Leftrightarrow (1 + r) \mathbf{m}^\top \boldsymbol{\lambda}_\Delta^{(0)} = 1 \\ &\Leftrightarrow (1 + r) \frac{\mathbf{m}^\top \exp \boldsymbol{\lambda}^{(0)}}{1 + \mathbf{m}^\top \exp \boldsymbol{\lambda}^{(0)}} = 1 \quad (27) \\ &\Leftrightarrow (1 + r) \mathbf{m}^\top \exp \boldsymbol{\lambda}^{(0)} = \\ &\quad 1 + \mathbf{m}^\top \exp \boldsymbol{\lambda}^{(0)} \\ &\Leftrightarrow \mathbf{m}^\top \exp \boldsymbol{\lambda}^{(0)} = \frac{1}{r} \end{aligned}$$

Assuming that $\boldsymbol{\lambda}^{(0)} = \omega \mathbf{1}_d$ with $\omega \in \mathbb{R}$, we get:

$$\begin{aligned} \mathbf{m}^\top \exp \boldsymbol{\lambda}^{(0)} &= \frac{1}{r} \Leftrightarrow \|\mathbf{m}\|_1 \exp \omega = \frac{1}{r} \\ &\Leftrightarrow \omega = -\log(r \|\mathbf{m}\|_1) \end{aligned} \quad (28)$$

Ablation To test the effect of these modifications on PDPGD and ALMA prox, we perform an ablation study on the masking strategy. We perform this experiment with DeepLabV3+ ResNet-50 on Cityscapes. We compare three different constraint masking strategies:

- masking only the unlabeled regions, denoted by \mathbf{m} ;
- masking the unlabeled regions and $(100 - \nu)\%$ of the largest constraints, denoted by $\tilde{\mathbf{m}}$;
- masking the unlabeled regions and linearly decreasing the fraction of constraints to reach $\nu\%$ at the last iteration, corresponding to strategy described in Equation (6), denoted by $\tilde{\mathbf{m}}^{(t)}$.

The results of this experiment are provided in Figure 5. While the effect is small for this particular model and dataset, the $\tilde{\mathbf{m}}$ and $\tilde{\mathbf{m}}^{(t)}$ strategies do result in smaller perturbations. However, PDPGD does not seem to benefit from the more advanced masking strategies. It obtains the best results when masking only the unlabeled regions. This issue comes from the projection of the dual variables onto the d -simplex (26): as different constraints get discarded in subsequent iterations, their relative weights vary drastically, leading to oscillations of the dual variables. This phenomenon does not occur with an Augmented Lagrangian based attacks, as the penalty multipliers are not projected together. This does not create a dependency between the multipliers, resulting in a more stable optimization. Therefore, for the experiments, ALMA prox is used with the adaptive masking strategy with a linear decay, whereas PDPGD is used with the unlabeled region masking only.

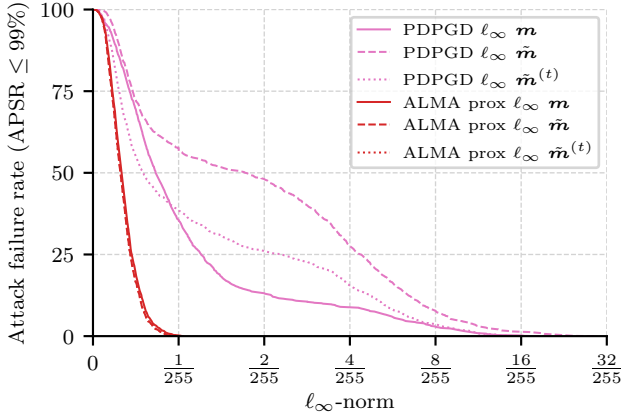


Figure 5. Influence of the constraint masking strategies on PDPGD and ALMA prox for untargeted attacks on Cityscapes with DeepLabV3+ ResNet-50. \mathbf{m} corresponds to masking the pixels with no labels, $\tilde{\mathbf{m}}$ corresponds to additionally masking the top $(100 - \nu)\%$ constraints, and $\tilde{\mathbf{m}}^{(t)}$ corresponds to the strategy in Equation (6). The curves for PDPGD and ALMA prox $\tilde{\mathbf{m}}$ and $\tilde{\mathbf{m}}^{(t)}$ overlap.

G. Attack complexities and run-times

Tables 4 and 5 report the average complexities (in terms of forward and backward) and run-time per sample for the attacks. The run-times for the targeted attacks are similar to their targeted variant, as the only difference lies in the loss, which is far from being the most computationally expensive part of the attacks. These results show that ALMA prox has slightly higher run-times compared to the other attacks with similar budgets, while being more effective. Note that some variance come from the fact that the experiments were run on a shared compute cluster.

Attack	Forwards	Backwards
I-FGSM [27] 13×20 and MI-FGSM [22]	260	260
PGD 13×40 and $13 \times 4 \times 10$ [29]	520	520
DAG $\eta = 0.003$ [39]	63	62
DAG $\eta = 0.001$ [39]	155	154
FMN ℓ_∞ [34]	500	500
PDPGD ℓ_∞ [30]	500	500
ALMA prox	500	500

Table 4. Average complexity of the attacks in terms of number of forward and backward propagations. All the attacks, except I-FGSM, MI-FGSM and PGD have a 500 iteration budget. DAG is the only attack that uses an early stopping criterion.

H. Complete attack results

Table 6 reports the median and average ℓ_∞ -norm (multiplied by 255 for readability) of the perturbations produced by the attacks for all regular models. As in Section 6, the perturbation norm is considered to be 1 for unsuccessful attacks. This means that attack with less than 50% success

have a 1 (*i.e.* 255) median ℓ_∞ -norm in the table. Additionally, Figures 6 and 7 show the percentage of unsuccessful attacks on Pascal VOC 2012 and Cityscapes for all models considered.

Attack	Pascal VOC 2012				Cityscapes					
	DeepLabV3+ ResNet-50	DeepLabV3+ ResNet-101	FCN HRNetV2 W48	DeepLabV3 DDC-AT	DeepLabV3+ ResNet-50	FCN HRNetV2 W48	SegFormer MiT-B0	SegFormer MiT-B3	DeepLabV3 DDC-AT	
Untargeted	I-FGSM 13×20 [27]	13.3	19.2	28.8	–	87.0	43.2	31.5	83.4	–
	MI-FGSM 13×20 [22]	13.3	19.2	29.1	–	87.6	41.5	31.1	83.2	–
	PGD CE 13×40 [29]	26.4	38.2	57.0	–	171.9	81.9	61.7	164.6	–
	PGD CE $13 \times 4 \times 10$ [29]	26.6	38.3	57.2	–	174.1	81.9	61.6	164.5	–
	PGD DLR 13×40 [29]	27.7	39.3	60.9	–	179.5	88.1	67.6	170.7	–
	PGD DLR $13 \times 4 \times 10$ [29]	27.9	39.4	57.6	–	179.3	88.2	67.6	170.6	–
	DAG $\eta = 0.003$ [39]	2.1	5.2	6.0	7.0	19.5	14.4	9.2	49.2	22.4
	DAG $\eta = 0.001$ [39]	4.8	12.9	14.8	12.7	45.5	37.3	23.4	107.3	44.7
	FMN ℓ_∞ [34]	24.4	32.0	41.3	24	155.2	79.1	59.1	158.2	64.5
	PDPGD ℓ_∞ [30]	28.1	35.5	41.0	25.5	161.0	84.5	64.1	162.4	64.7
	ALMA prox	28.1	36.1	43.4	28	161.9	96.8	77.6	176.3	69.8
Targeted	I-FGSM 13×20 [27]	13.3	19.2	29.8	–	87.5	41.4	31.1	83.0	–
	MI-FGSM 13×20 [22]	13.3	19.4	28.8	–	87.8	41.4	31.1	83.1	–
	PGD CE 13×40 [29]	27.0	38.1	56.8	–	171.9	81.9	61.4	165.0	–
	PGD CE $13 \times 4 \times 10$ [29]	26.7	38.4	56.5	–	171.9	82.0	61.5	164.3	–
	PGD DLR 13×40 [29]	27.7	39.4	58.9	–	179.6	88.2	68.4	172.0	–
	PGD DLR $13 \times 4 \times 10$ [29]	27.8	39.8	58.7	–	179.4	88.2	69.7	171.0	–
	DAG $\eta = 0.003$ [39]	0.8	1.5	1.9	2.2	47.6	18.6	21.9	54.1	44.5
	DAG $\eta = 0.001$ [39]	2.1	4.2	4.6	4.8	111.7	44.9	52.0	132.9	61.0
	FMN ℓ_∞ [34]	24.4	32.1	41.1	25.4	155.1	79.2	58.7	157.9	60.4
	PDPGD ℓ_∞ [30]	27.6	36.9	39.0	26.8	158.8	81.4	64.1	159.1	63.2
	ALMA prox	28.4	36.3	43.0	28.5	161.9	97.2	78.4	179.5	67.6

Table 5. Average run-times per image for the attacks, in seconds.

Attack	Pascal VOC 2012						Cityscapes								
	DeepLabV3+ ResNet-50		DeepLabV3+ ResNet-101		FCN HRNetV2 W48		DeepLabV3+ ResNet-50		FCN HRNetV2 W48		SegFormer MiT-B0		SegFormer MiT-B3		
Untargeted	I-FGSM 13×20 [27]	146.32	136.55	124.14	131.34	227.11	145.69	255.00	242.15	255.00	194.98	106.34	121.12	255.00	208.97
	MI-FGSM 13×20 [22]	195.35	145.96	188.37	150.13	255.00	157.63	255.00	244.43	255.00	218.91	202.15	159.52	255.00	228.43
	PGD CE 13×40 [29]	80.10	120.93	100.42	123.17	152.08	136.05	255.00	236.77	255.00	176.94	9.31	48.10	255.00	188.15
	PGD CE $13 \times 4 \times 10$ [29]	24.90	74.15	20.05	30.28	66.93	118.76	255.00	244.92	255.00	216.38	34.77	42.49	255.00	231.51
	PGD DLR 13×40 [29]	7.22	26.42	9.00	21.29	11.11	23.85	17.75	24.23	13.29	16.15	23.53	32.34	84.22	102.80
	PGD DLR $13 \times 4 \times 10$ [29]	4.42	24.99	6.41	11.02	10.46	29.75	255.00	227.89	12.08	45.17	22.37	35.49	110.82	134.41
	DAG $\eta = 0.003$ [39]	5.69	6.63	6.65	8.74	8.80	10.61	6.30	7.51	8.42	9.91	6.41	6.48	8.83	9.02
	DAG $\eta = 0.001$ [39]	5.23	8.22	6.17	21.14	8.49	14.61	5.95	9.39	8.20	20.29	6.08	13.57	8.59	53.65
	FMN ℓ_∞ [34]	0.46	38.34	0.56	51.60	0.91	46.39	1.97	96.57	0.97	35.12	0.58	1.16	1.08	6.42
	PDPGD ℓ_∞ [30]	0.73	1.77	1.17	3.49	1.52	2.43	1.06	2.82	1.69	2.75	0.87	1.17	1.39	3.05
	ALMA prox	0.32	0.34	0.37	0.41	0.51	0.56	0.24	0.26	0.40	0.41	0.26	0.26	0.33	0.33
Targeted	I-FGSM 13×20 [27]	0.47	0.50	0.65	0.67	0.59	0.64	255.00	255.00	115.76	128.65	5.91	31.96	51.73	88.24
	MI-FGSM 13×20 [22]	0.59	0.66	0.77	0.82	0.77	0.85	4.26	55.35	4.23	4.47	2.48	2.53	2.54	2.60
	PGD CE 13×40 [29]	0.37	0.43	0.50	0.55	0.50	0.54	3.49	3.71	3.49	3.62	2.30	2.34	1.87	1.92
	PGD CE $13 \times 4 \times 10$ [29]	0.50	0.51	0.62	0.70	0.62	0.65	255.00	255.00	255.00	255.00	255.00	255.00	255.00	255.00
	PGD DLR 13×40 [29]	0.62	0.68	0.81	0.92	0.81	0.94	8.37	67.86	7.41	7.52	4.79	4.96	3.98	4.09
	PGD DLR $13 \times 4 \times 10$ [29]	0.87	1.07	1.18	2.03	1.12	1.28	255.00	255.00	255.00	255.00	255.00	255.00	255.00	255.00
	DAG $\eta = 0.003$ [39]	4.21	4.50	4.84	5.09	5.32	5.66	11.34	12.96	9.87	10.49	8.05	8.44	9.82	10.06
	DAG $\eta = 0.001$ [39]	3.92	4.21	4.55	5.80	5.07	5.36	10.96	40.28	9.43	14.42	255.00	145.61	11.53	119.47
	FMN ℓ_∞ [34]	0.42	0.45	0.46	0.56	0.47	0.49	255.00	254.36	2.25	2.34	255.00	255.00	255.00	255.00
	PDPGD ℓ_∞ [30]	0.28	0.36	0.31	0.40	0.35	0.46	14.51	14.41	16.71	16.75	17.93	17.87	19.26	19.20
	ALMA prox	0.25	0.26	0.29	0.30	0.32	0.34	1.15	1.17	1.11	1.12	0.70	0.70	0.65	0.66

Table 6. Median and average $\|\delta\|_\infty \times 255$ for each adversarial attack on Pascal VOC 2012 and Cityscapes for the regular models.

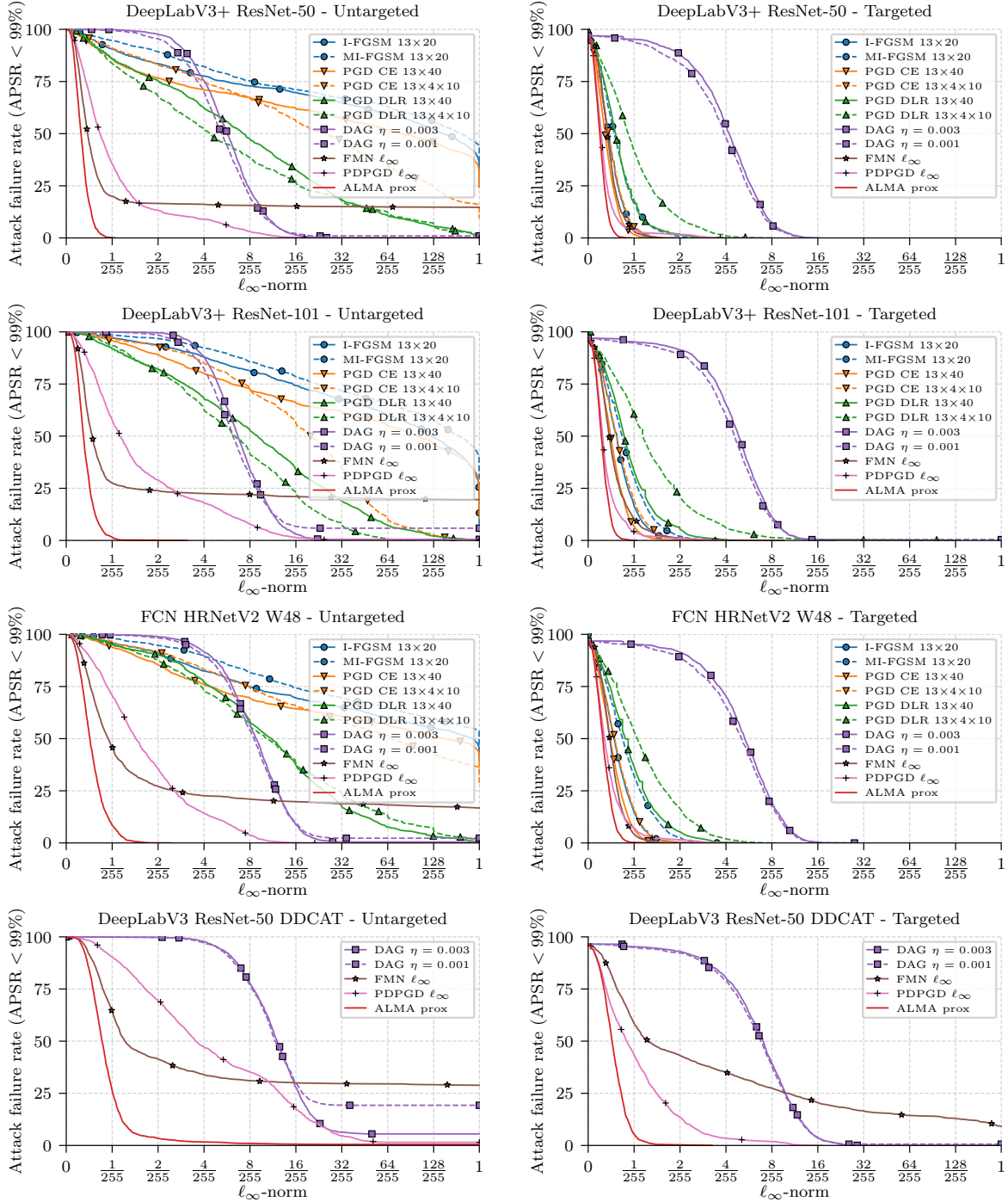


Figure 6. Percentage of unsuccessful ℓ_∞ attacks on **Pascal VOC 2012** (i.e. with $\text{APSR} \leq 99\%$). A stronger attack has a lower curve; a more robust model has a higher curve. Horizontal axis is linear on $[0, 2/255]$ and logarithmic on $[2/255, 1]$.

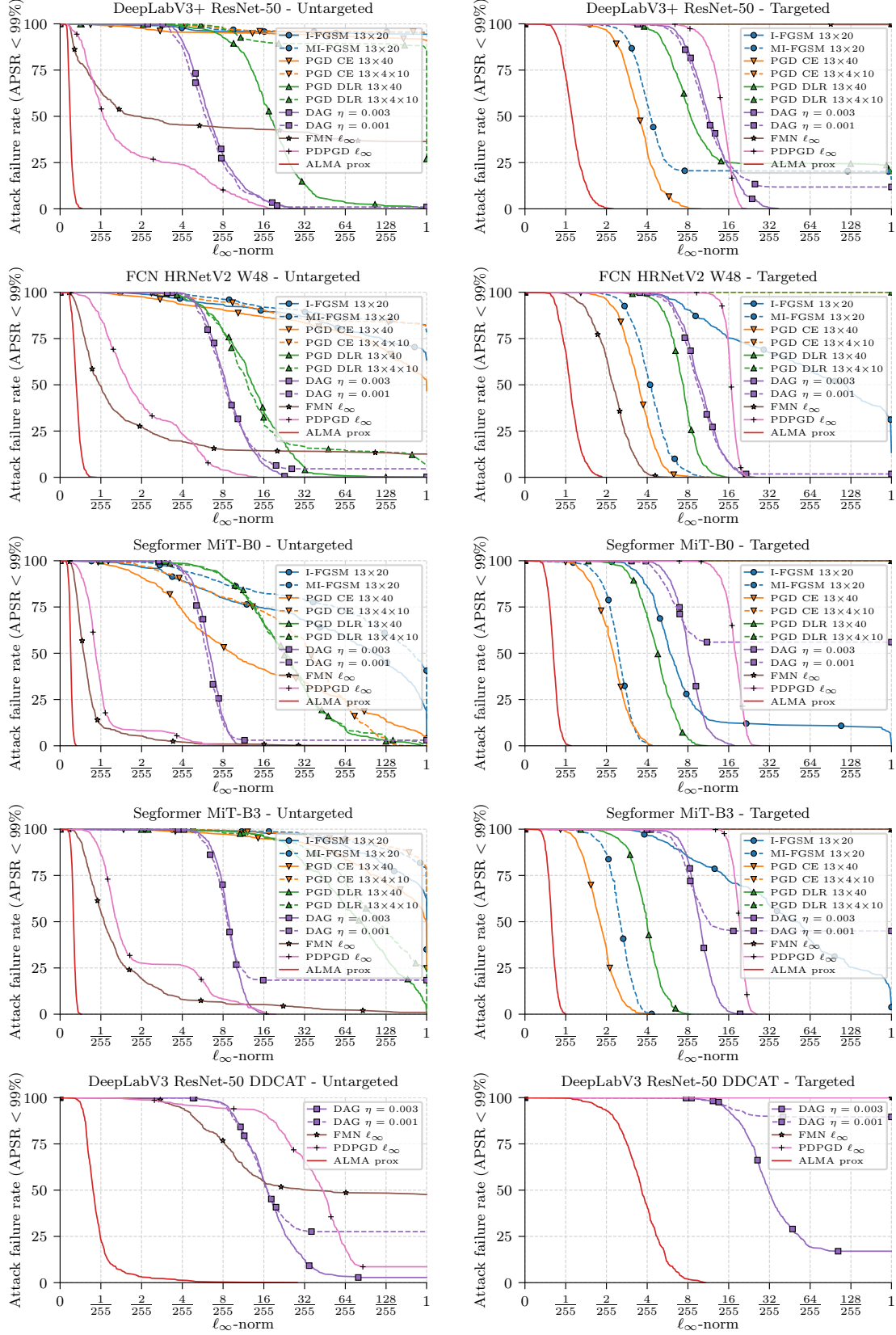


Figure 7. Percentage of unsuccessful ℓ_∞ attacks on **Cityscapes** (*i.e.* with $\text{APSR} \leq 99\%$). A stronger attack has a lower curve; a more robust model has a higher curve. Horizontal axis is linear on $[0, 2/255]$ and logarithmic on $[2/255, 1]$.