

Theoretic Fundamentals of Machine and Deep Learning

Certified Robustness II: Ablations on Randomized Smoothing. High Dimensions and Computer Vision

Aleksandr Petiushko

Lomonosov MSU, Faculty of Mechanics and Mathematics
MIPT, RAIKI
Nuro, Autonomy Interaction Research

Winter-Spring, 2023



Content

- ① Certified robustness: recap
- ② Certification in High Dimensional case
- ③ Certification of Semantic Perturbations

AP

Certified Robustness: for classification

- Let us NN function $f(x)$ is the classifier to K classes: $f : \mathbb{R}^d \rightarrow Y, Y = \{1, \dots, K\}$
- Usually we have NN $h(x) : \mathbb{R}^d \rightarrow \mathbb{R}^K$, and $f(x) = \arg \max_{i \in Y} h(x)_i$

Deterministic approach

Need to find the class of input perturbation $S(x, f)$ so as the classifier's output doesn't not change, or more formally:

$$f(x + \delta) = f(x) \quad \forall \delta \in S(x, f)$$

Probabilistic approach

Need to find the class of input perturbation $S(x, f, P)$ w.r.t. robustness probability P s.t.:

$$\text{Prob}_{\delta \in S(x, f, P)}(f(x + \delta) = f(x)) = P$$

Remark: Probabilistic approach coincides with Deterministic one when $P = 1$.

Certified Robustness: for regression

- Let us NN function $f(x)$ NN $f(x)$ is the regressor: $f : \mathbb{R}^d \rightarrow \mathbb{R}$

Deterministic approach

Need to find the class of input perturbation $S(x, f, f_{low}, f_{up})$ w.r.t. the upper and lower bounds on the output perturbation f_{low}, f_{up} s.t.:

$$f(x) - f_{low} \leq f(x + \delta) \leq f(x) + f_{up} \quad \forall \delta \in S(x, f, f_{low}, f_{up})$$

Probabilistic approach

Need to find the class of input perturbation $S(x, f, f_{low}, f_{up}, P)$ w.r.t. robustness probability P and the upper / lower bounds on the output perturbation f_{low}, f_{up} s.t.:

$$\text{Prob}_{\delta \in S(x, f, f_{low}, f_{up}, P)}(f(x) - f_{low} \leq f(x + \delta) \leq f(x) + f_{up}) = P$$

Certified Robustness: inverse tasks for classification

- Suppose that we know the input perturbation class S
- For classification we have only probabilistic formulation

Classification

Need to measure the probability P of retaining the classifier's output under some class of input perturbations S :

$$\text{Prob}_{\delta \in S}(f(x + \delta) = f(x)) = P$$

Certified Robustness: inverse tasks for regression

- Suppose that we know the input perturbation class S
- For regression we have both deterministic and probabilistic formulations

Regression (deterministic formulation)

Need to find the upper and lower bounds $f_{low}(f, x, S), f_{up}(f, x, S)$ of the output perturbation under some class of input perturbations S :

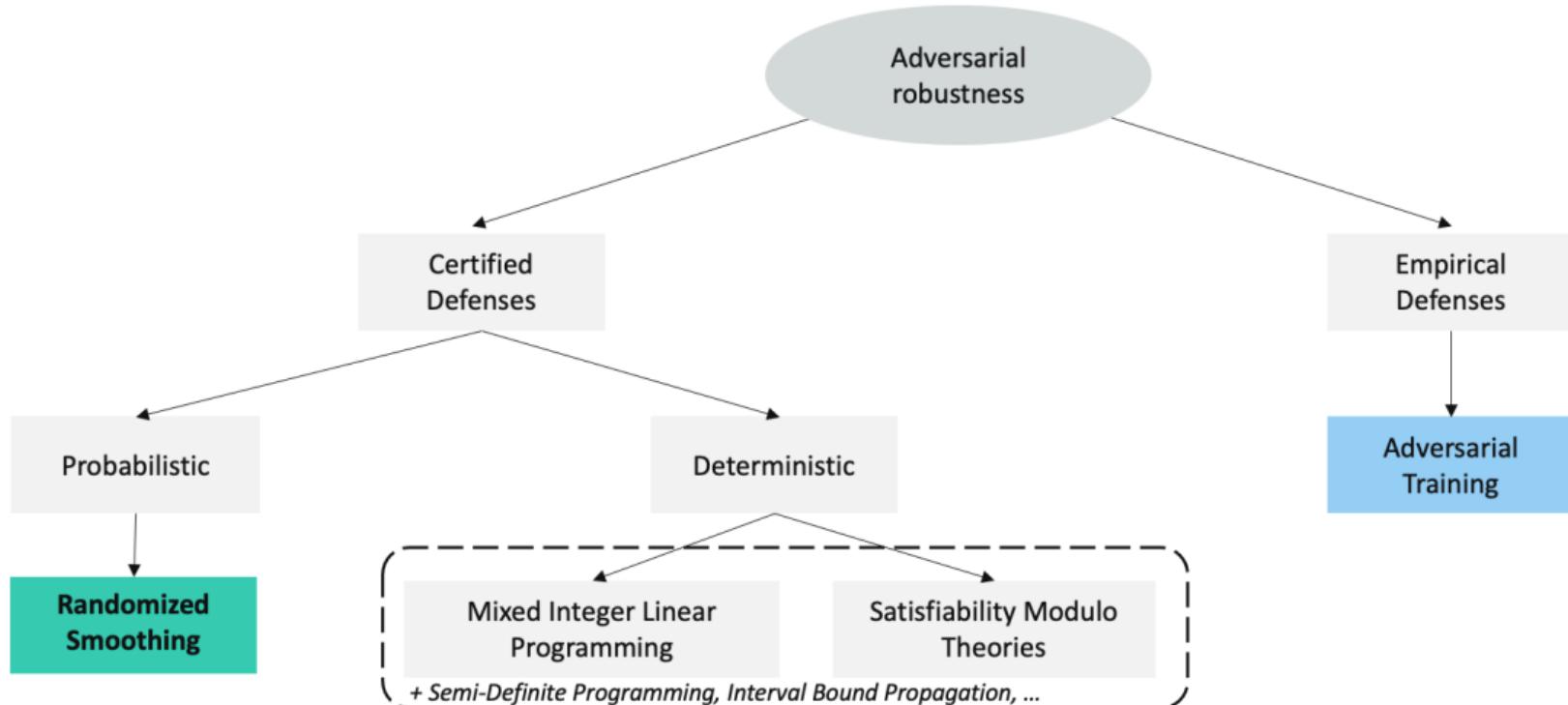
$$f(x) - f_{low}(f, x, S) \leq f(x + \delta) \leq f(x) + f_{up}(f, x, S)$$

Regression (probabilistic formulation)

Need to measure the probability P of keeping the classifier's output inside the lower / upper bounds f_{low}, f_{up} under some class of input perturbations S :

$$\text{Prob}_{\delta \in S} (f(x) - f_{low} \leq f(x + \delta) \leq f(x) + f_{up}) = P$$

Adversarial Robustness: overview



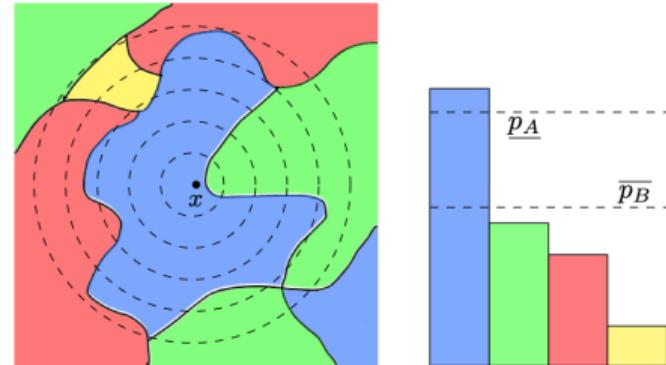
Applicable only for *very small NN models* – e.g. for MNIST/CIFAR

Randomized Smoothing

Idea of Randomized Smoothing (RS)

- Let's use the **Test Time Augmentation (TTA)** in order to mitigate the boundary effect
- The new classifier $g(x)$ is defined as:

$$g(x) = \arg \max_{c \in Y} P(f(x + \epsilon) = c), \epsilon \sim N(0, \sigma^2)$$



RS main result

- If the initial classifier $f(x)$ is robust under Gaussian noise,
- Then the new classifier $g(x)$ is robust under **ANY** noise

Randomized Smoothing: Theory overview

Theorem: Certification Radius

Suppose $c_A \in Y$ and $\underline{p}_A, \overline{p}_B \in [0, 1]$ satisfy

$\mathbb{P}(f(x + \epsilon) = c_A) \geq \underline{p}_A \geq \overline{p}_B \geq \max_{c_B \neq c_A} \mathbb{P}(f(x + \epsilon) = c_B)$. Then
 $g(x + \delta) = c_A \quad \forall \|\delta\|_2 < R$, where

$$R = \frac{\sigma}{2} (\Phi^{-1}(\underline{p}_A) - \Phi^{-1}(\overline{p}_B))$$

Remark. Φ^{-1} is the inverse of the standard Gaussian CDF: $\Phi(x) = \frac{1}{2\pi} \int_{-\infty}^x e^{-\frac{t^2}{2}} dt$.

Randomized Smoothing: Better training

- For an original RS the training process is just augmentation with Gaussian Noise images
- Idea of **SmoothAdv**¹: let's do adversarial training (AT) using attacks on smoothed classifier $g(x)$!
 - ▶ Original RS hard example (in the vicinity ϵ):

$$x' = \arg \max_{x': \|x' - x\|_2 \leq \epsilon} \mathbb{E}_{\delta \sim N(0, \sigma^2 I)} [L(f_\theta(x' + \delta), y)]$$

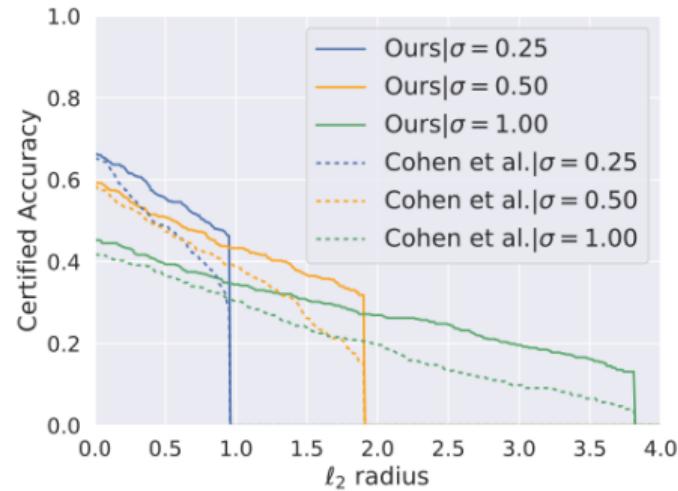
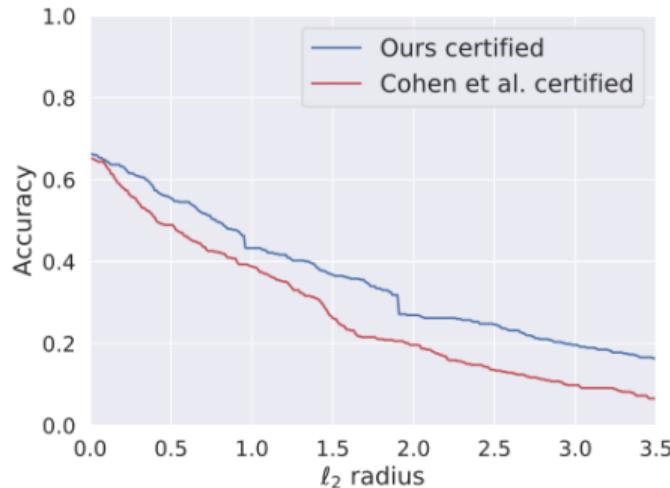
- ▶ SmoothAdv hard example:

$$x' = \arg \max_{x': \|x' - x\|_2 \leq \epsilon} L(g_\theta(x'), y) = \arg \max_{x': \|x' - x\|_2 \leq \epsilon} L(\mathbb{E}_{\delta \sim N(0, \sigma^2 I)} [f_\theta(x' + \delta)], y)$$

- ▶ If cross entropy loss is used, then the difference is that we changing $\sum \log$ to $\log \sum \Rightarrow$ using Jensen's inequality, we have $\sum \log \leq \log \sum$

¹Salman, Hadi, et al. "Provably robust deep learning via adversarially trained smoothed classifiers." 2019.

Randomized Smoothing: SmoothAdv Results

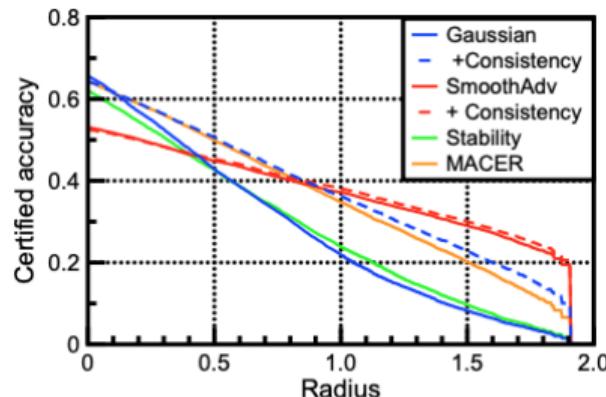


ℓ_2 RADIUS (IMAGENET)	0.5	1.0	1.5	2.0	2.5	3.0	3.5
COHEN ET AL. [6] (%)	49	37	29	19	15	12	9
OURS (%)	56	45	38	28	26	20	17

Randomized Smoothing: Regularization

- What if we can work on top of smoothed classifier $g(x)$ to make it more reasonable?
- Idea of Consistency Regularization²: let's force *similarity* between *smoothed* and *perturbed* predictions as well as minimizing the entropy of smoothed output:

$$L_{CR}(x) = \lambda \mathbb{E}_{\delta \sim N(0, \sigma^2 I)} D_{KL}(g(x) || f(x + \delta)) + \eta H(g(x))$$



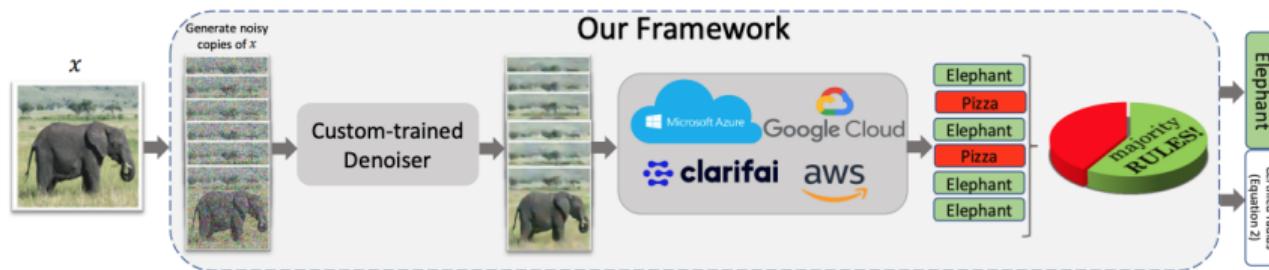
(b) $\sigma = 0.50$

²Jeong, Jongheon, and Jinwoo Shin. "Consistency regularization for certified robustness of smoothed classifiers." 2020

Randomized Smoothing: Black-box access

- What if we **cannot change the pretrained classifier**, but want to increase its certified robustness?
- Idea of **Black-box smoothing**³: Let's train a **denoiser** D used after we've added Gaussian noise!
 - ▶ And then simply apply the majority rule

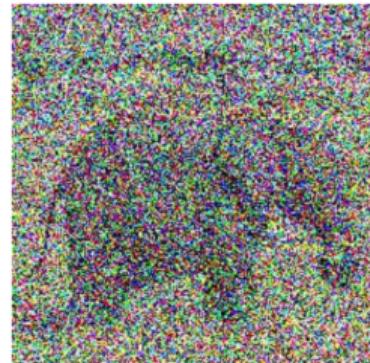
$$g(x) = \arg \max_{c \in Y} \mathbb{P}[f(D(x + \delta)) = c], \quad \delta \sim N(0, \sigma^2 I)$$



³Salman, Hadi, et al. “Black-box smoothing: A provable defense for pretrained classifiers.” 2020

Randomized Smoothing: Denoiser for Black-box

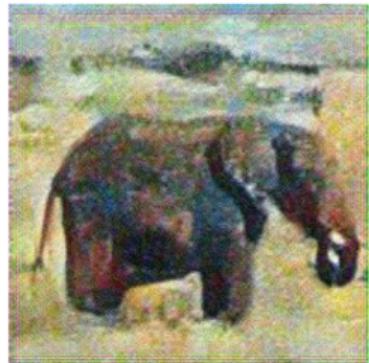
- Denoiser: trained with two losses for every Gaussian σ :
 - ▶ MSE
 - ▶ Stability (classification cross entropy)



(a) Noisy



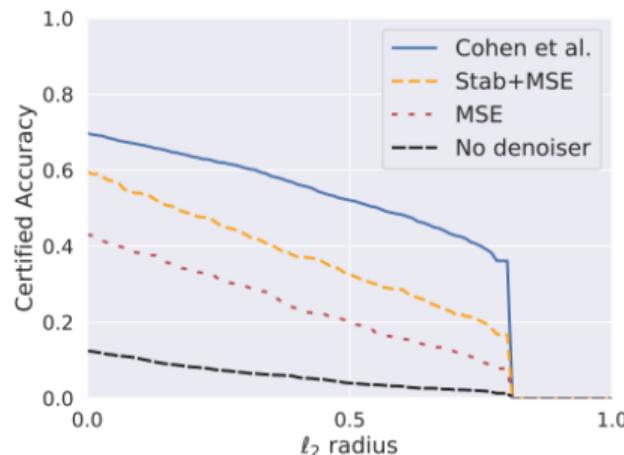
(b) MSE



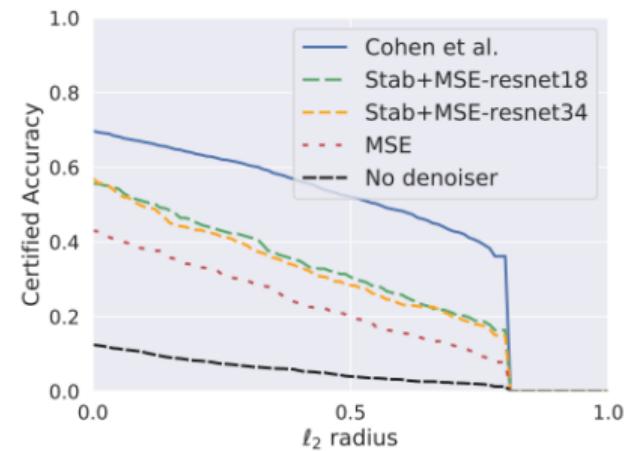
(c) Stab+MSE

Randomized Smoothing: Results for Black-box

Full access to classifier during training



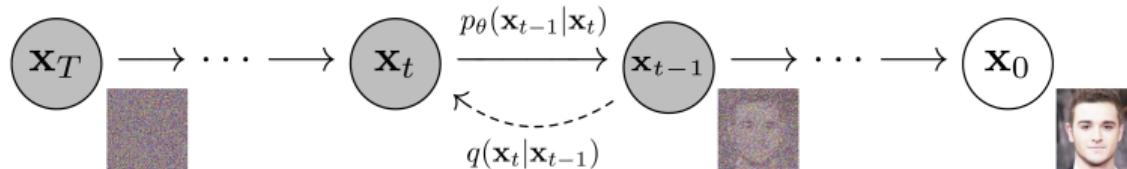
Using surrogate classifier during training



ℓ_2 RADIUS (IMAGENET)	0.25	0.5	0.75	1.0	1.25	1.5
COHEN ET AL. (2019) (%)	(70) 62	(70) 52	(62) 45	(62) 39	(62) 34	(50) 29
NO DENOISER (BASELINE) (%)	(49) 32	(12) 4	(12) 2	(0) 0	(0) 0	(0) 0
OURS (BLACK-BOX) (%)	(69) 48	(56) 31	(56) 19	(34) 12	(34) 7	(30) 4
OURS (WHITE-BOX) (%)	(67) 50	(60) 33	(60) 20	(38) 14	(38) 11	(38) 6

Denoiser by DDPM⁵

- A novel approach⁴ to use off-the-shelf models:
 - ▶ SotA classifier (trained on clean images)
 - ▶ Denoising Diffusion Model
 - ★ Based on the noise level σ , estimate $\bar{\alpha}_t, t$
 - ★ Generate $x_t \sim N(\sqrt{\bar{\alpha}_t} \cdot x, (1 - \bar{\alpha}_t)I)$
 - ★ Denoise by DDPM decoder (using **only 1 step**): $\hat{x} = \text{denoise}(x_t)$
 - ★ Classify!
- Results in 14% improvement over the prior certified SoTA, and an improvement of 30% over denoised smoothing



⁴N. Carlini, F. Tramer, and Z. Kolter. "(Certified!!) Adversarial Robustness for Free!", 2022

⁵J. Ho, A. Jain, and P. Abbeel. "Denoising diffusion probabilistic models", 2020

Randomized Smoothing: vector functions

- Previously all results were for the classifiers: $f, g : \mathbb{R}^d \rightarrow Y, Y = \{1, \dots, K\}$,
 $g(x) = \arg \max_{c \in Y} P(f(x + \epsilon) = c), \epsilon \sim N(0, \sigma^2)$
- Let's consider the vector-based functions f (e.g., feature vector): $\mathbb{R}^d \rightarrow \mathbb{R}^D$
- Then the smoothed version g of it we'll define as: $g(x) = \mathbb{E}_{\epsilon \sim N(0, \sigma^2 I)}[f(x + \epsilon)]$
- In this case the following relation to Lipschitz functions can be established⁶:

Lipschitz-continuity of smoothed vector function

Suppose that $g(x)$ is continuously differentiable for all x . If for all x , $\|f(x)\|_2 = 1$, then $g(x)$ is L -Lipschitz in l_2 -norm with $L = \sqrt{\frac{2}{\pi\sigma^2}}$.

Randomized Smoothing: adversarial embedding risk

- Let's establish the beautiful geometrical fact useful for the few-shot classification:

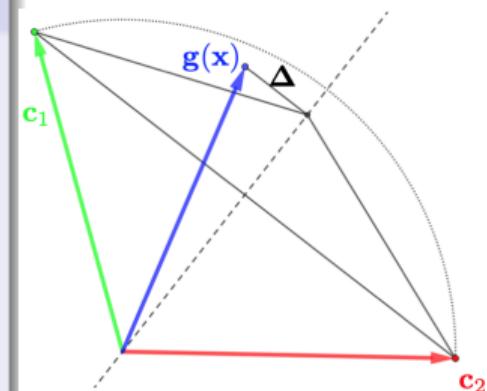
Adversarial embedding risk

Given an input $x \in \mathbb{R}^d$ and the embedding $g : \mathbb{R}^d \rightarrow \mathbb{R}^D$ the closest point on to decision boundary in the embedding space is located at a distance:

$$\gamma = \|\Delta\|_2 = \frac{\|c_2 - g(x)\|_2^2 - \|c_1 - g(x)\|_2^2}{2\|c_2 - c_1\|_2^2},$$

where $c_1 \in \mathbb{R}^D$ and $c_2 \in \mathbb{R}^D$ are the two closest prototypes.

- γ is the distance between classifying embedding and the decision boundary between classes represented by c_1 and c_2 .
- γ is the minimum l_2 -distortion in the embedding space required to change the prediction of g .



Randomized Smoothing: certification

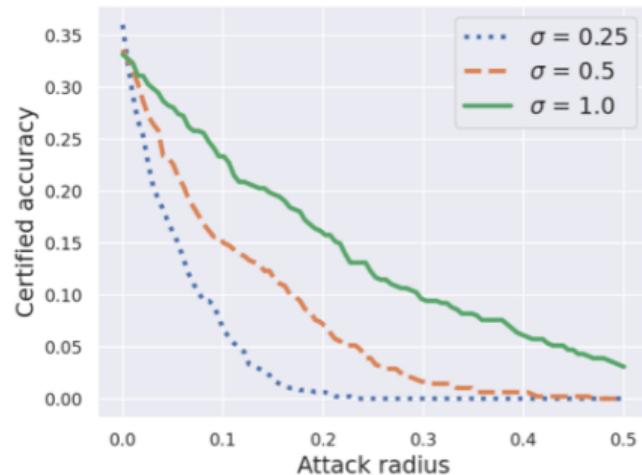
- Two results above lead to the certification guarantee:

Robustness guarantee

Certified radius r of g at x , where g is the smoothed version of $f : \|f(x)\|_2 = 1$, is

$$r = \frac{\gamma}{L}$$

1-shot results for *miniImageNet*⁷



⁷Vinyals, Oriol, et al. "Matching networks for one shot learning." 2016

Randomized Smoothing: Some Results on Regression

- Let function $f : \mathbb{R}^d \rightarrow \mathbb{R}$
- Smoothed version $g : \mathbb{R}^d \rightarrow \mathbb{R}$, $g(x) = \mathbb{E}_{\epsilon \sim N(0, \sigma^2)}[f(x + \epsilon)]$
- Φ^{-1} is the inverse of the standard Gaussian CDF: $\Phi(x) = \frac{1}{2\pi} \int_{-\infty}^x e^{-\frac{t^2}{2}} dt$
- Result: the solution⁸ for inverse certification problem (deterministic formulation):

Theorem

For any bounded $f : \mathbb{R}^d \rightarrow [l, u]$, the map $\eta(x) = \sigma\Phi^{-1}(\frac{g(x)-l}{u-l})$ is 1-Lipschitz, implying

$$l + (u - l) \cdot \Phi\left(\frac{\eta(x) - \|\delta\|_2}{\sigma}\right) \leq g(x + \delta) \leq l + (u - l) \cdot \Phi\left(\frac{\eta(x) + \|\delta\|_2}{\sigma}\right)$$

Exercise. Prove it.

⁸Chiang, Ping-yeh, et al. "Detection as regression: Certified object detection with median smoothing." AP
2020

Randomized Smoothing: norms

- Randomized Smoothing = Smoothing distribution + **norm** l_p of perturbation
- Using l_p -balls is neither necessary nor sufficient for perceptual robustness
- Certification is only for much smaller regions than humans can do
- Remark about physical nature of l_p -balls:
 - ▶ l_2 corresponds to the power of signals
 - ▶ l_1 corresponds to the pixel mass
 - ▶ l_∞ corresponds to the noise in camera sensors
 - ▶ l_0 corresponds to the practical patch robustness

Randomized Smoothing: High Dimensional Case

- The perturbation δ is measured by l_p -norm
- $p = 1$ and $p = 2$ are the only **special cases**⁹: $R = \frac{\sigma}{2}(\Phi^{-1}(\underline{p}_A) - \Phi^{-1}(\overline{p}_B))$
- Unfortunately, these are **only** examples of **non-decreasing** with **input dimension** d
- For any $p \geq 2$, the certification radius¹⁰ is decreasing with dimensionality d :

$$R_p(x) = \frac{\sigma}{2d^{\frac{1}{2} - \frac{1}{p}}}(\Phi^{-1}(\underline{p}_A) - \Phi^{-1}(\overline{p}_B))$$

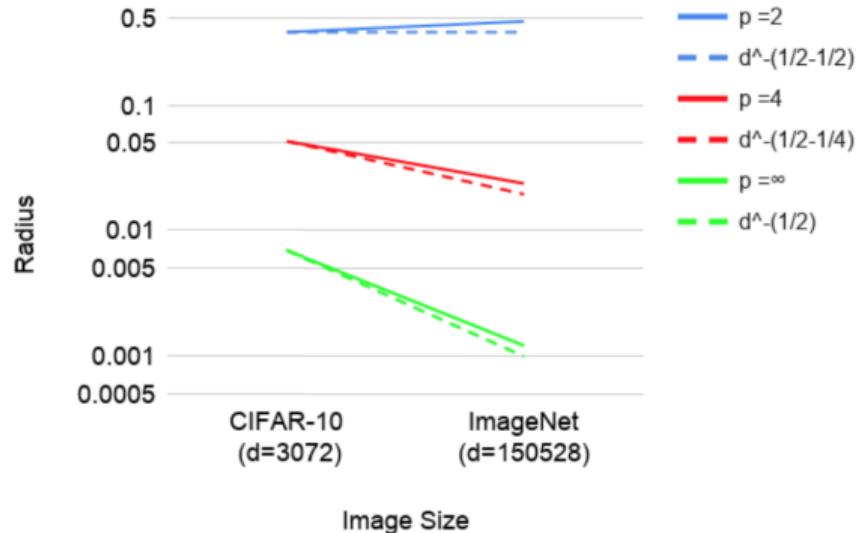
- And the most important case for Computer Vision (CV), $p = \infty$, means

$$R_\infty \sim \frac{1}{\sqrt{d}}$$

⁹Yang, Greg, et al. "Randomized smoothing of all shapes and sizes." 2020

¹⁰Kumar, Aounon, et al. "Curse of dimensionality on randomized smoothing for certifiable robustness." AP
2020

Randomized Smoothing: CV illustration



High Dimension Case in CV

- Any **semantic-meaningful** perturbation in **CV** leads to **high l_∞ -perturbation**, and the dimension of an image $d = H \times W$ usually is very high (like millions of pixels)
- $R_\infty \sim \frac{1}{\sqrt{d}}$ means that there is **no any practical certified radius**
 - ▶ E.g., it is hard to achieve promising robust accuracy ($\geq 70\%$) even when the perturbation radius is as small as 2 pixels¹¹
- E.g., for semantic-specific transformations like **contrast** and **brightness** the **error is higher** than on clean images up to 50-60% on *Common Corruptions*¹² on ImageNet



Network	Error	Bright	Contrast
AlexNet	43.5	100	100
SqueezeNet	41.8	97	98
VGG-11	31.0	75	86
VGG-19	27.6	68	80
VGG-19+BN	25.8	61	74
ResNet-18	30.2	69	78
ResNet-50	23.9	57	71

¹¹Blum, Avrim, et al. "Random smoothing might be unable to certify l_∞ robustness for high-dimensional images." 2020

¹²Hendrycks, Dan, and Thomas Dietterich. "Benchmarking neural network robustness to common corruptions and perturbations." 2019

High Dimension Case in CV: Autonomous Driving

- The same is true for safety-critical applications like autonomous driving¹³



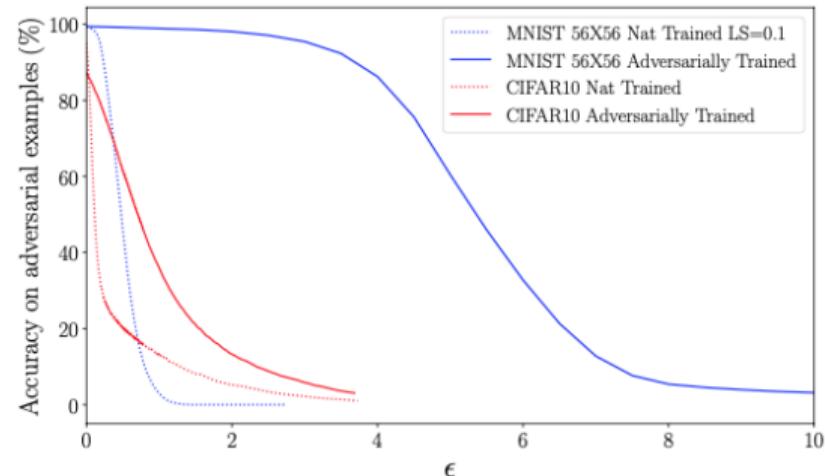
Transformation	#err
Brightness	97
Contrast	31

¹³Tian, Yuchi, et al. "Deeptest: Automated testing of deep-neural-network-driven autonomous cars."

High Dimension Case: dataset complexity¹⁴

- Let $C \sim \frac{1}{var}$, where var is dataset variance
- Then certified radius ϵ exists:
 - ▶ On $(n - 1)$ -dim sphere with probability $C \cdot e^{-\frac{n-1}{2}\epsilon^2}$
 - ▶ On n -dim hypercube with probability $C \cdot \frac{e^{-2\pi\epsilon^2}}{2\pi\epsilon}$ for $p \geq 2$
 - ▶ On l_0 -ball with probability $C \cdot e^{-\frac{\epsilon^2}{n}}$
- **Conclusion1:** Highly concentrated datasets (with big C / small var) can be relatively safe from adversarial examples

Conclusion2: Data distribution, and not dimensionality, is the primary cause of adversarial susceptibility



¹⁴Shafahi, Ali, et al. "Are adversarial examples inevitable?" 2018

Semantic perturbations for additive parameters

- So... let's certify semantic perturbations¹⁵!
 - ▶ Usually parameterized by a much smaller dimension (1 or 2 dimensional)
- Consider **rotations** and **translations** γ_β parameterized by β : $\gamma_\beta : \mathbb{R}^d \rightarrow \mathbb{R}^d$
- A smoothed classifier $g(x) = \arg \max_{c \in Y} P_{\beta \sim N(0, \sigma^2)}(f \circ \gamma_\beta(x) = c)$
- Also **interpolation** procedure is taken into account because after rotation we need to interpolate anyway

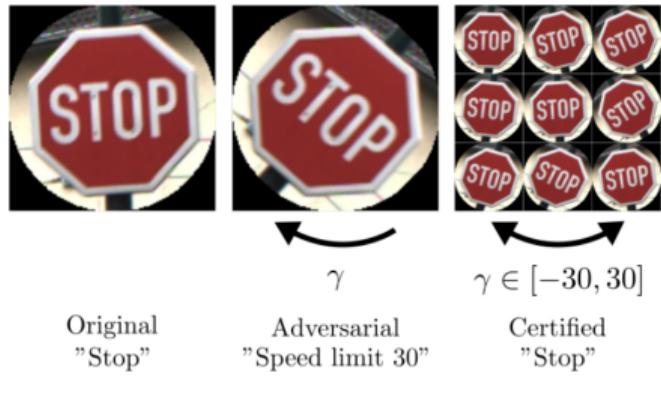
Certification Radius

Suppose $c_A \in Y$ and $\underline{p}_A, \overline{p}_B \in [0, 1]$ satisfy

$\mathbb{P}_{\beta \sim N(0, \sigma^2)}(f \circ \gamma_\beta(x) = c_A) \geq \underline{p}_A \geq \overline{p}_B \geq \max_{c_B \neq c_A} \mathbb{P}_{\beta \sim N(0, \sigma^2)}(f \circ \gamma_\beta(x) = c_B)$. Then
 $g \circ \gamma_\beta(x) = c_A \quad \forall \|\beta\|_2 < r_\gamma$, where $r_\gamma = \frac{\sigma}{2}(\Phi^{-1}(\underline{p}_A) - \Phi^{-1}(\overline{p}_B))$.

¹⁵Fischer, Marc, et al. "Certified defense to image transformations via randomized smoothing." 2020   

Semantic perturbations for additive parameters: results



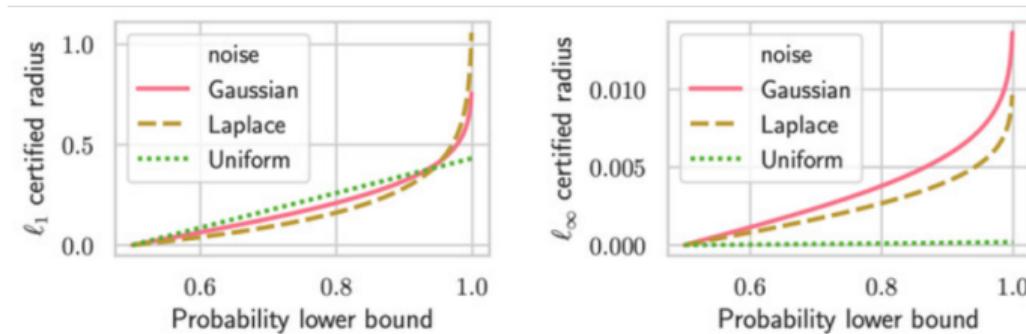
Dataset	\mathcal{I}	σ_γ	α_γ	Rotation		r_γ percentile		
				<i>f</i> Acc.	<i>g</i> Acc.	25 th	50 th	75 th
ImageNet	bil.	10	0.001	0.39	0.29	10.81	10.81	10.81
ImageNet	bil.	10	0.001	0.39	0.29	18.29	18.29	18.29
ImageNet	bil.	30	0.001	0.39	0.28	9.09	16.59	28.60
ImageNet	bil.	30	0.001	0.39	0.28	20.22	25.36	30 [†]
ImageNet	bic.	10	0.001	0.39	0.29	10.40	10.40	10.40
ImageNet	bic.	30	0.001	0.39	0.27	9.33	17.00	28.74
ImageNet	near.	10	0.001	0.39	0.29	9.62	9.62	9.62
ImageNet	near.	30	0.001	0.39	0.26	7.38	16.63	27.72

Dataset	\mathcal{I}	σ_γ	α_γ	Translation		r_γ percentile		
				<i>f</i> Acc.	<i>g</i> Acc.	25 th	50 th	75 th
ImageNet	bil.	50	0.001	0.48	0.36	2.4%	2.4%	2.4%
ImageNet	bic.	50	0.001	0.48	0.36	2.4%	2.4%	2.4%

AP

Randomized Smoothing: smoothing distribution

- Randomized Smoothing = Smoothing **distribution** + norm l_p of perturbation
- Original (and most of the follow-up ones) work uses Gaussian Smoothing
- Other types of smoothing could be taking into account: e.g. Uniform¹⁶ or Laplacian¹⁷
 - ▶ E.g., when doing uniform smoothing inside l_∞ or l_1 -ball of radius b , it is proved¹⁸
 $R_p < \frac{2b}{d^{1-\frac{1}{p}}}$ or $R_p < \frac{2b}{d}$ correspondingly
- What about other types?



¹⁶Lee, Guang-He, et al. "Tight certificates of adversarial robustness for randomly smoothed classifiers."

2019

¹⁷Teng, Jiaye, et al. " ℓ_1 Adversarial Robustness Certificates: a Randomized Smoothing Approach." 2019 AP

¹⁸Kumar, A., et al. "Curse of dimensionality on randomized smoothing for certifiable robustness." 2020

Semantic perturbations and multiplicative parameters

- All research above is concentrated on **additive** perturbations
- Let's investigate the **multiplicative** parameters¹⁹ (e.g., *gamma correction* $G_\gamma(x) = x^\gamma$ in CV)
- **Definition:** A parameterized map $\psi_\delta : X \rightarrow X$, $\delta \in \mathcal{B} \subset \mathbb{R}^n$ is called multiplicatively composable if $(\psi_\delta \circ \psi_\theta)(x) = \psi_{(\delta \cdot \theta)}(x)$, $\forall x \in X$, $\forall \delta, \theta \in \mathcal{B}$
- Example: $G_\beta \circ G_\gamma(x) = (x^\gamma)^\beta = x^{\gamma \cdot \beta} = G_{\gamma \cdot \beta}(x)$



¹⁹Muravev, Nikita, and Aleksandr Petiushko. "Certified Robustness via Randomized Smoothing over Multiplicative Parameters." 2021

Semantic perturbations and multiplicative parameters: results

- To work under this limitation, the new type of smoothing distribution is needed:
 - Positive support
 - Mean at 1
- The proposal to use is **Rayleigh** distribution:
 $p_\beta(z) = \sigma^{-2}ze^{-z^2/(2\sigma^2)}, z \geq 0$
- Then the following is true: $g \circ \psi_\gamma(x) = c_A$ for all γ satisfying $\gamma_1 < \gamma < \gamma_2$, where γ_1, γ_2 are the only solutions of the following equations:
 $F(\gamma_1^{-1}F^{-1}(\overline{p_B})) + F(\gamma_1^{-1}F^{-1}(1 - \underline{p_A})) = 1,$
 $F(\gamma_2^{-1}F^{-1}(\underline{p_A})) + F(\gamma_2^{-1}F^{-1}(1 - \overline{p_B})) = 1,$
and $F(z) = 1 - e^{-z^2/(2\sigma^2)}$ is the CDF of γ .
- The results are better for $\gamma < 1$ in comparison to Uniform, Gaussian and Laplace smoothing

$\underline{p_A}$	$\overline{p_B}$	γ_1	γ_2
0.600	0.400	0.86	1.15
	0.200	0.71	1.33
0.700	0.300	0.72	1.32
	0.100	0.54	1.56
0.800	0.200	0.57	1.52
0.900	0.100	0.39	1.82
0.990	0.010	0.12	2.58
0.999	0.001	0.04	3.16

Semantic perturbations and compositions

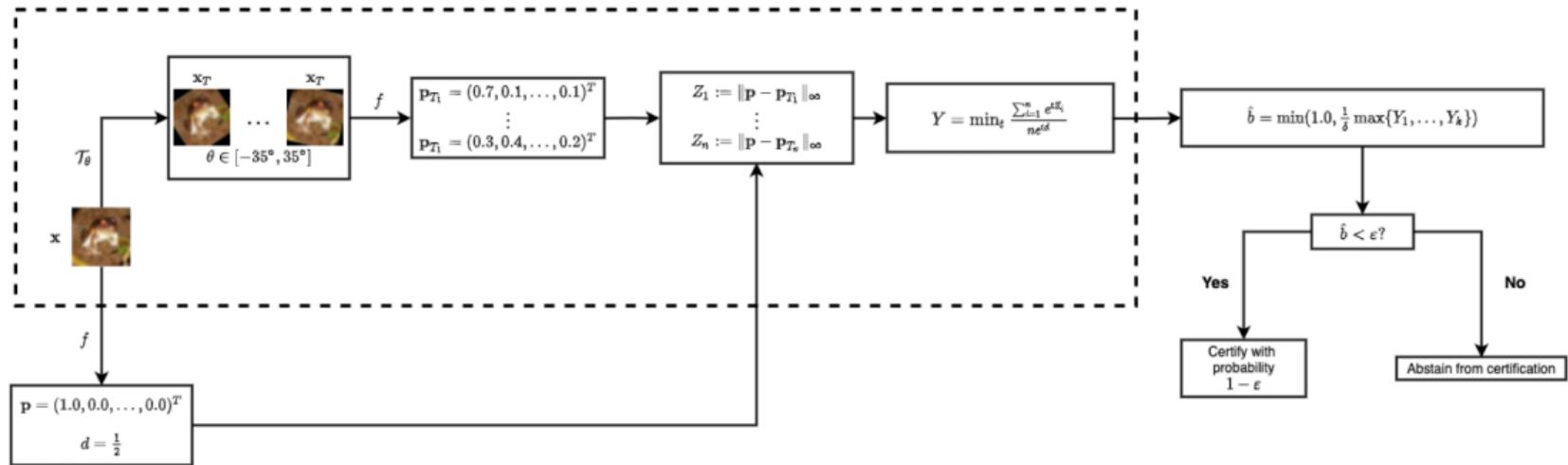
- Usually multiple transformations are applied to the input: how to certify the composition?
- Forward theoretical estimation is difficult \Rightarrow let's try inverse (probabilistic) task²⁰!
- The proposal to use Chernoff-Cramer inequality²¹ (Markov's inequality corollary) to provide the statistically-grounded **estimations for the certification**, where perturbed radius is already given
- Can be easily used for **any semantic perturbation** and **any compositions**

Dataset	Transform	Parameters	Training type	ERA	PCA(ε)		
					$\varepsilon = 10^{-10}$	$\varepsilon = 10^{-7}$	$\varepsilon = 10^{-4}$
	Brightness	$\theta_b \in [-40\%, 40\%]$	plain	58.4%	47.8%	51.6%	55.2%
				smoothing	65.0%	55.4%	59.4%
	Contrast	$\theta_c \in [-40\%, 40\%]$	plain	91.6%	62.4%	67.0%	69.6%
				smoothing	88.0%	67.0%	72.8%
	Rotation	$\theta_r \in [-10^\circ, 10^\circ]$	plain	73.4%	64.6%	69.0%	71.0%
				smoothing	72.4%	57.4%	63.6%
Contrast + Brightness	see Contrast & Brightness		plain	0.0%	0.0%	0.0%	0.0%
				smoothing	0.4%	0.0%	0.0%
	Rotation + Brightness	see Rotation & Brightness	plain	22.6%	16.2%	20.6%	21.8%
				smoothing	30.4%	21.2%	24.6%
Scale + Brightness	see Scale & Brightness		plain	10.2%	10.4%	10.4%	10.4%
				smoothing	41.8%	40.6%	40.6%

²⁰Pautov, Mikhail, et al. "CC-Cert: A probabilistic approach to certify general robustness of neural networks." 2021

²¹Boucheron, Stéphane, et al. "Concentration inequalities." 2003

Inverse certification for any transformation²²



²²Pautov, Mikhail, et al. "CC-Cert: A probabilistic approach to certify general robustness of neural networks." 2021

Semantic perturbations: further development

- Later works introduced approaches to take into account different types of perturbations and interpolation errors²³

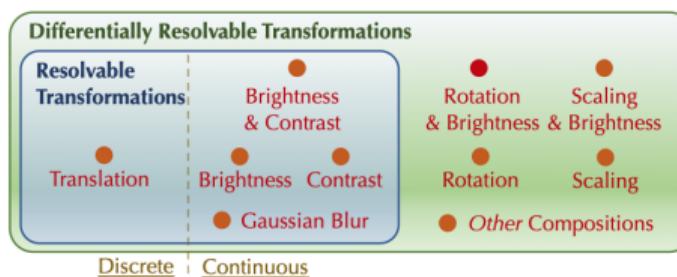
DEFINITION 2 (RESOLVABLE TRANSFORM). A transformation $\phi: \mathcal{X} \times \mathcal{Z} \rightarrow \mathcal{X}$ is called resolvable if for any $\alpha \in \mathcal{Z}$ there exists a resolving function $\gamma_\alpha: \mathcal{Z} \rightarrow \mathcal{Z}$ that is injective, continuously differentiable, has non-vanishing Jacobian and for which

$$\phi(\phi(x, \alpha), \beta) = \phi(x, \gamma_\alpha(\beta)) \quad x \in \mathcal{X}, \beta \in \mathcal{Z}. \quad (7)$$

Furthermore, we say that ϕ is additive, if $\gamma_\alpha(\beta) = \alpha + \beta$.

DEFINITION 3 (DIFFERENTIALLY RESOLVABLE TRANSFORM). Let $\phi: \mathcal{X} \times \mathcal{Z}_\phi \rightarrow \mathcal{X}$ be a transformation with noise space \mathcal{Z}_ϕ and let $\psi: \mathcal{X} \times \mathcal{Z}_\psi \rightarrow \mathcal{X}$ be a resolvable transformation with noise space \mathcal{Z}_ψ . We say that ϕ can be resolved by ψ if for any $x \in \mathcal{X}$ there exists function $\delta_x: \mathcal{Z}_\phi \times \mathcal{Z}_\phi \rightarrow \mathcal{Z}_\psi$ such that for any $\alpha \in \mathcal{Z}_\phi$ and any $\beta \in \mathcal{Z}_\phi$,

$$\phi(x, \alpha) = \psi(\phi(x, \beta), \delta_x(\alpha, \beta)). \quad (15)$$



²³Li, Linyi, et al. "Tss: Transformation-specific smoothing for robustness certification." 2020

Randomized Smoothing: Object Detection²⁵

- The approach treats certification for **Object Detection** as a **Regression** problem for **Black-box**²⁴ detectors
- In order to certify multiple BB under different noise, the **sorting** based on location + **binning** based on label is needed
- But the **certification** is still **out of practical use**

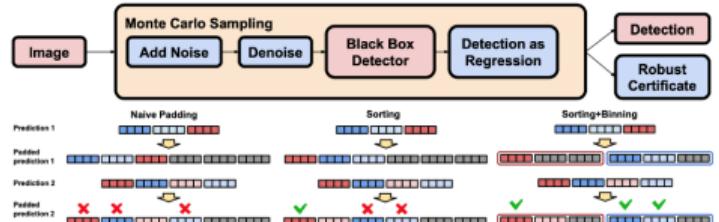
Architecture	Base Detector		Smoothed Detector	
	AP @ 50	AP @ 50	Certified AP @ 50	Certified AP @ 50
YOLOv3	48.66%	31.93%	4.21%	4.21%
Mask RCNN	51.28%	30.53%	1.67%	1.67%
Faster RCNN	50.47%	29.89%	1.54%	1.54%

²⁴Salman, H., et al. "Black-box smoothing: A provable defense for pretrained classifiers." 2020

²⁵Chiang, P., et al. "Detection as regression: Certified object detection with median smoothing." 2020 ↗ ↙ ↘

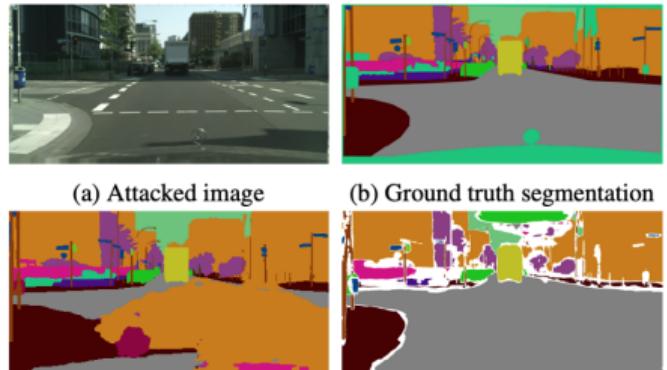


Figure 1: Samples of object detection certificates using the proposed method. Dotted lines represent the farthest a bounding box could move under an adversarial perturbation δ of bounded ℓ_2 -norm. If the predicted bounding box can be made to disappear, or if the label can be changed, after a perturbation with $\|\delta\|_2 < 0.36$, then we annotate the bounding box with a red X.



Randomized Smoothing: Semantic Segmentation²⁶

- For **segmentation every pixel** needs to be classified. Any atomic error leads to the **overall certification fail**
- Authors proposed to **allow some pixels** to be **non-classified** if the smoothed classifier provides less probability than $\tau \in (\frac{1}{2}]$ and redefine it:
 - $\bar{f}_i^\tau = c_{A,i}$ if $\mathbb{P}_{\epsilon \sim N(0,\sigma^2)}(f_i(x + \epsilon)) > \tau$ and \emptyset else.
 - $c_{A,i} = \arg \max_{c \in Y} \mathbb{P}_{\epsilon \sim N(0,\sigma^2)}(f_i(x + \epsilon) = c)$
- Making this assumption, they succeed to prove the certification Theorem and provide non-trivial results on segmentation benchmarks:
 - Theorem.** Let $I_x = \{i | \bar{f}_i^\tau \neq \emptyset, i = 1, \dots, N\}$ denote the set of non-abstain indices for \bar{f}^τ . Then, $\bar{f}_i^\tau(x + \delta) = \bar{f}_i^\tau(x) \forall x \in I_x$ for $\delta \in \mathbb{R}^{N \times 2}$ with $\|\delta\|_2 \leq R = \sigma \Phi^{-1}(\tau)$

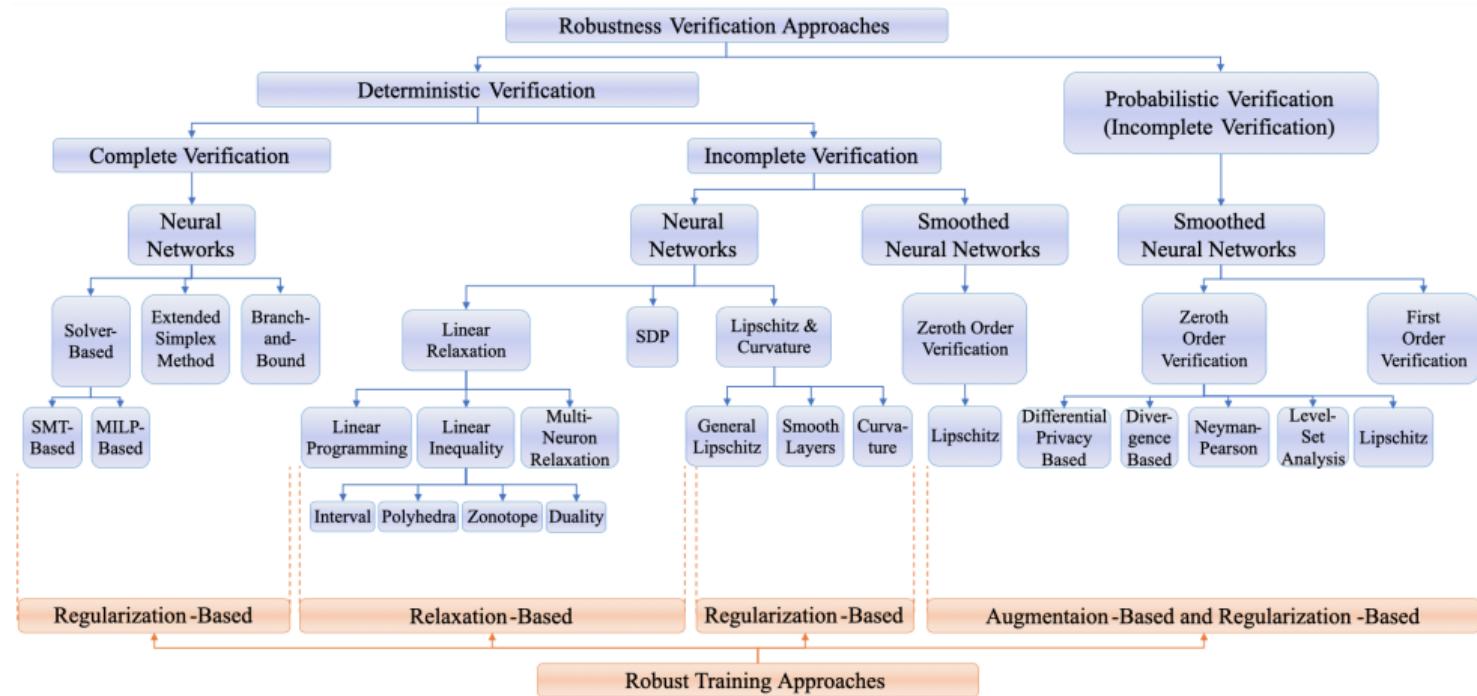


(c) Attacked segmentation (d) Certified segmentation

scale	σ	R	Cityscapes			
			acc.	mIoU	$\% \varnothing$	t
0.25	non-robust model	-	0.93	0.60	0.00	0.38
	base model	-	0.87	0.42	0.00	0.37
$n = 100, \tau = 0.75$	0.25	0.17	0.84	0.43	0.07	70.00
	0.33	0.22	0.84	0.44	0.09	70.21
	0.50	0.34	0.82	0.43	0.13	71.45
$n = 500, \tau = 0.95$	0.25	0.41	0.83	0.42	0.11	229.37
	0.33	0.52	0.83	0.42	0.12	230.69
	0.50	0.82	0.77	0.38	0.20	230.09

²⁶Fischer, Marc, et al. "Scalable certified segmentation via randomized smoothing." 2021

Systematization of Knowledge²⁷



²⁷SoK, Benchmark and Leaderboard

Takeaway notes

- Straightforward certification in l_∞ is not working for high dimension input
- In Computer Vision no need in any l_p (aside from l_0 for patch attacks, but it is usually also combined with other perturbations)
- Semantic perturbations are much harder to certify (+ interpolation!)
- **Current challenge:** 3D and even non-rigid transformations of **real world**

Thank you!