

OPTIMAL TRANSPORT AS A DEFENSE AGAINST
ADVERSARIAL ATTACKS

Quentin Bouniot, Romaric Audigier, Angélique Loesch

Université Paris-Saclay, CEA, List, F-91120, Palaiseau, France

Context

- **Adversarial attacks** : *human-imperceptible* perturbation for a given image to mislead a model.
- Most effective defenses based on adversarial training align *original* and *adversarial* representations.

Problems

- Defenses are *partially* aligning moments of distributions.
 - Can we **fully** align the distributions ?
- Current evaluation use a *fixed* perturbation size ϵ that can *differ* between papers.
 - How can we **choose** this perturbation size ?

Adversarial Examples

Inputs:

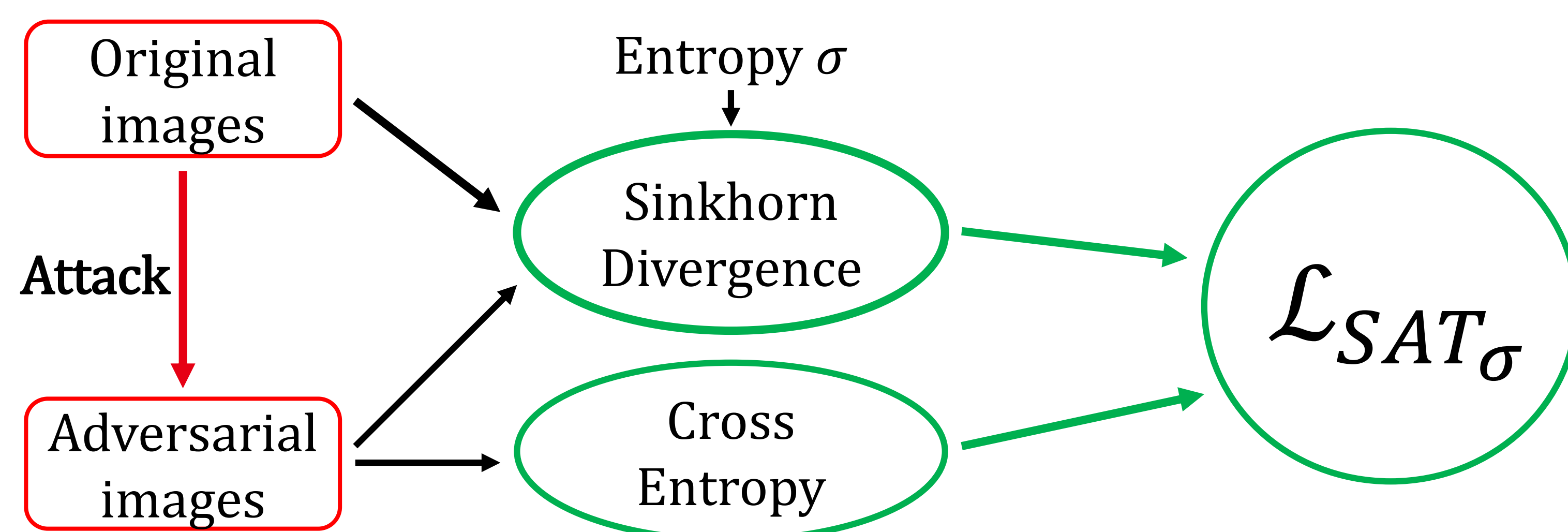
Original

 $\epsilon = 16$ Predictions:

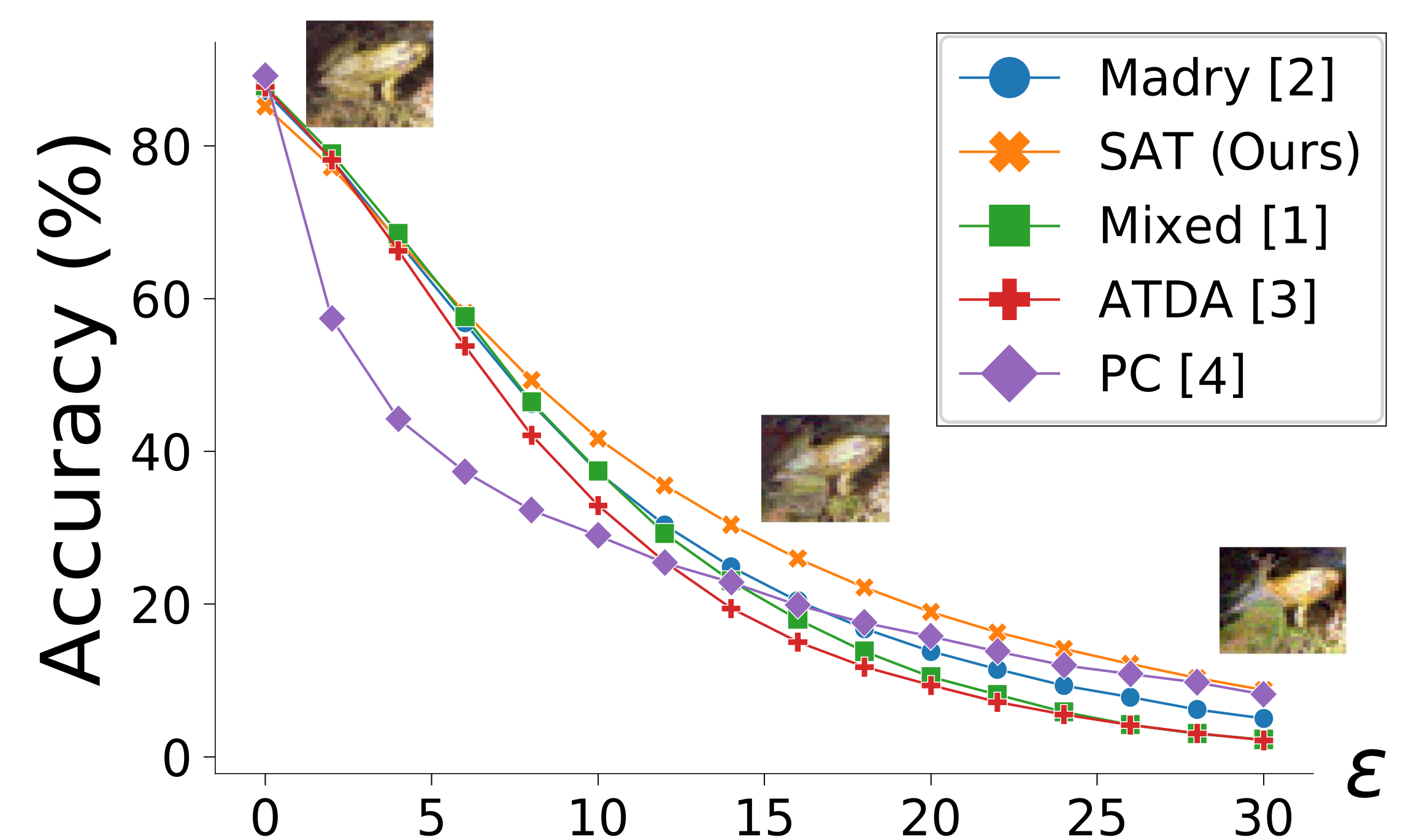
frog ✓

deer ✗

Sinkhorn Adversarial Training (SAT): a new Adversarial Defense



- Our **Sinkhorn Adversarial Training (SAT)** is based on theory of **Optimal Transport** [5] to consider the *whole* distributions and reflect *geometric properties*.



- A *fixed* perturbation size does not fully compare robustness.
- Our **SAT** is globally more robust than other SOTA defenses.

Area Under Accuracy Curve (AUAC): a new Metric for Robustness

- We propose **Area Under Accuracy Curve (AUAC)**:

$$AUAC_{\epsilon_{max}}(f) = \frac{1}{\epsilon_{max}} \int_{\epsilon=0}^{\epsilon_{max}} Acc(f, \epsilon, \mathbf{D}^{ts}) d\epsilon$$

- $Acc(f, \epsilon, \mathbf{D}^{ts})$ is the accuracy of f on the test set \mathbf{D}^{ts} with perturbations of size up to ϵ .
- AUAC quantifies more completely the robustness to adversarial attacks over a *wide range of perturbation sizes*.
- Evaluation also depends on the *Adversarial Attack* considered (see our paper for more examples).

Dataset	Archi.	Model	AUAC (%)	
			$\epsilon_{max} = 16$	$\epsilon_{max} = 30$
CIFAR-10	Resnet20	Non-defended	5.79	3.09
		Madry [2]	44.18	26.53
		Mixed [1]	40.68	22.73
		ATDA [3]	35.58	21.63
		SAT (Ours)	44.26	29.69
	Resnet110	PC [4]	37.89	26.47
CIFAR-100	WideResnet28-10	Non-defended	8.8	4.69
		Madry [2]	49.37	31.54
		Mixed [1]	49.27	30.01
		ATDA [3]	46.19	27.94
		SAT (Ours)	51.93	35.12
	WideResnet28-10	Non-defended	6.03	3.22
CIFAR-100	WideResnet28-10	Madry [2]	27.27	16.14
		Mixed [1]	27.80	16.13
		ATDA [3]	28.59	17.11
		SAT (Ours)	29.69	19.83

References

- [1] I. J. Goodfellow, J. Shlens, and C. Szegedy, "Explaining and harnessing adversarial examples," in International Conference on Learning Representations (ICLR), 2014.
- [2] A. Madry, A. Makelov, L. Schmidt, D. Tsipras, and A. Vladu, "Towards deep learning models resistant to adversarial attacks," in International Conference on Learning Representations (ICLR), 2018.
- [3] C. Song, K. He, L. Wang, and J. E. Hopcroft, "Improving the generalization of adversarial training with domain adaptation," in International Conference on Learning Representations (ICLR), 2019.

- [4] A. Mustafa, S. Khan, M. Hayat, R. Goecke, J. Shen, and L. Shao, "Adversarial defense by restricting the hidden space of deep neural networks," in International Conference on Computer Vision (ICCV), 2019.
- [5] J. Feydy, T. Séjourné, F.-X. Vialard, S.-i. Amari, A. Trounev, and G. Peyré, "Interpolating Between Optimal Transport and MMD using Sinkhorn Divergences," in Proceedings of Machine Learning Research (PMLR), 2019.