

Adversarial Examples and Adversarial Training

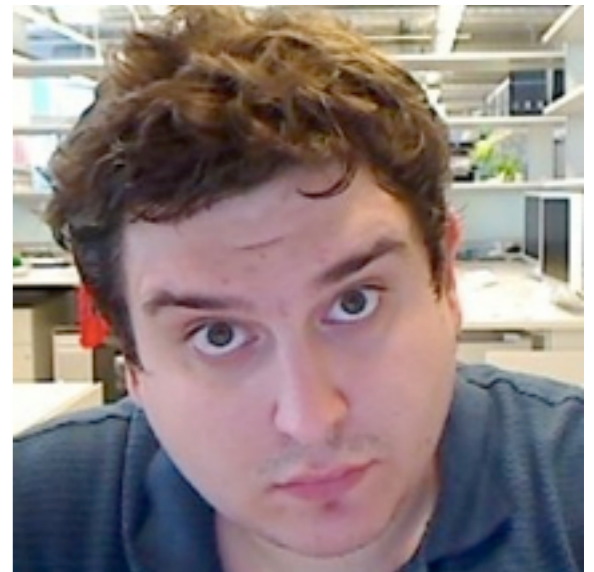
Ian Goodfellow, OpenAI Research Scientist

Presentation at San Francisco AI Meetup, 2016-08-18

OpenAI

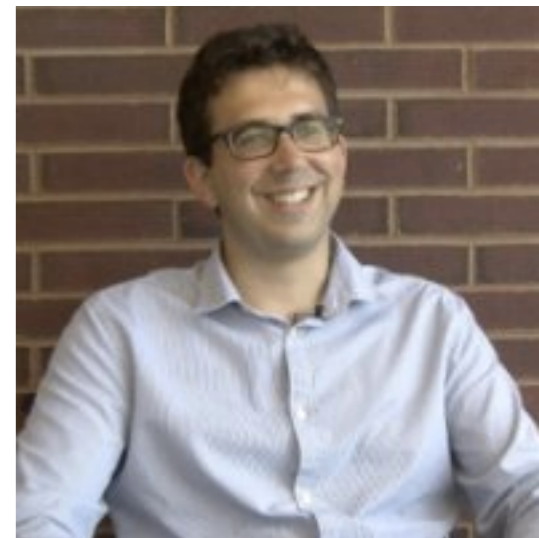
In this presentation

- “Intriguing Properties of Neural Networks” Szegedy et al, 2013
- “Explaining and Harnessing Adversarial Examples” Goodfellow et al 2014
- “Adversarial Perturbations of Deep Neural Networks” Warde-Farley and Goodfellow, 2016



In this presentation

- “Transferability in Machine Learning: from Phenomena to Black-Box Attacks using Adversarial Samples”
Papernot et al 2016
- “Practical Black-Box Attacks against Deep Learning Systems using Adversarial Examples” Papernot et al 2016
- “Adversarial Perturbations Against Deep Neural Networks for Malware Classification” Grosse et al 2016
(not my own work)



In this presentation

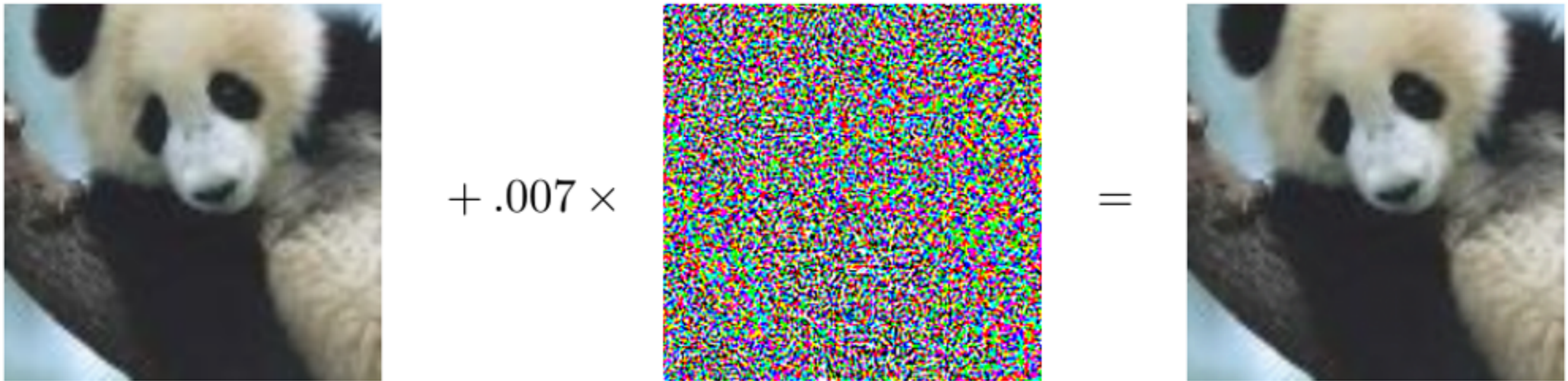
- “Distributional Smoothing with Virtual Adversarial Training” Miyato et al 2015 (**not my own work**)
- “Virtual Adversarial Training for Semi-Supervised Text Classification” Miyato et al 2016
- “Adversarial Examples in the Physical World” Kurakin et al 2016



Overview

- What are adversarial examples?
- Why do they happen?
- How can they be used to compromise machine learning systems?
- What are the defenses?
- How to use adversarial examples to improve machine learning, even when there is no adversary

Adversarial Examples



Timeline:

“Adversarial Classification” Dalvi et al 2004: fool spam filter

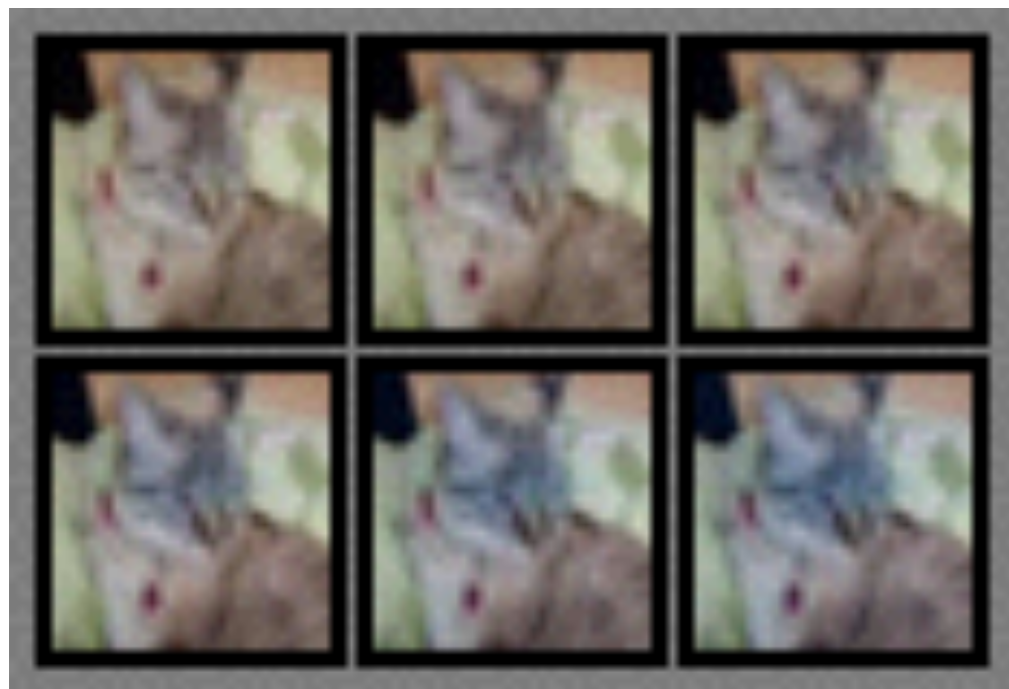
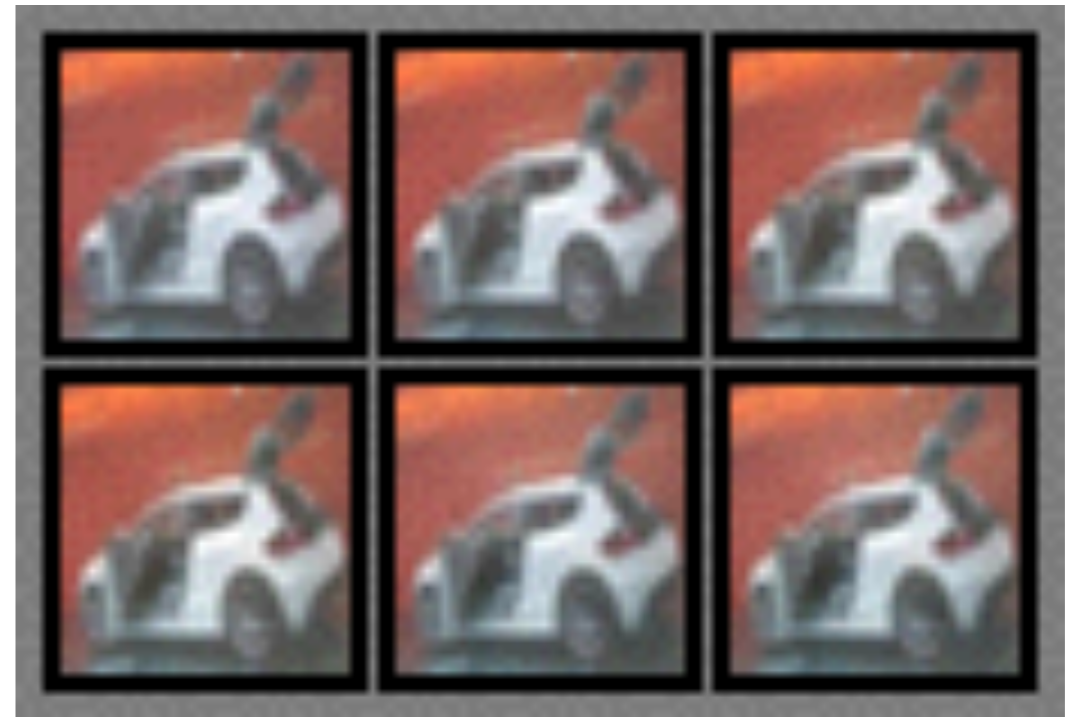
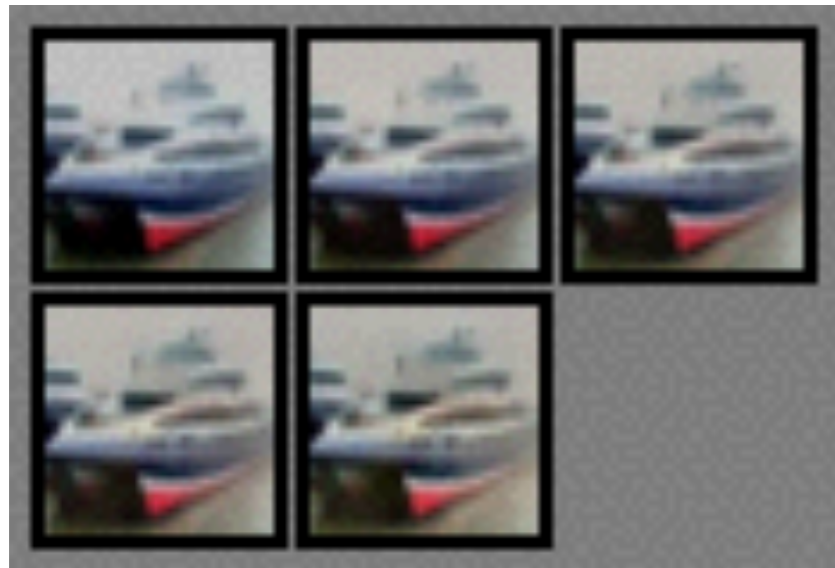
“Evasion Attacks Against Machine Learning at Test Time”

Biggio 2013: fool neural nets

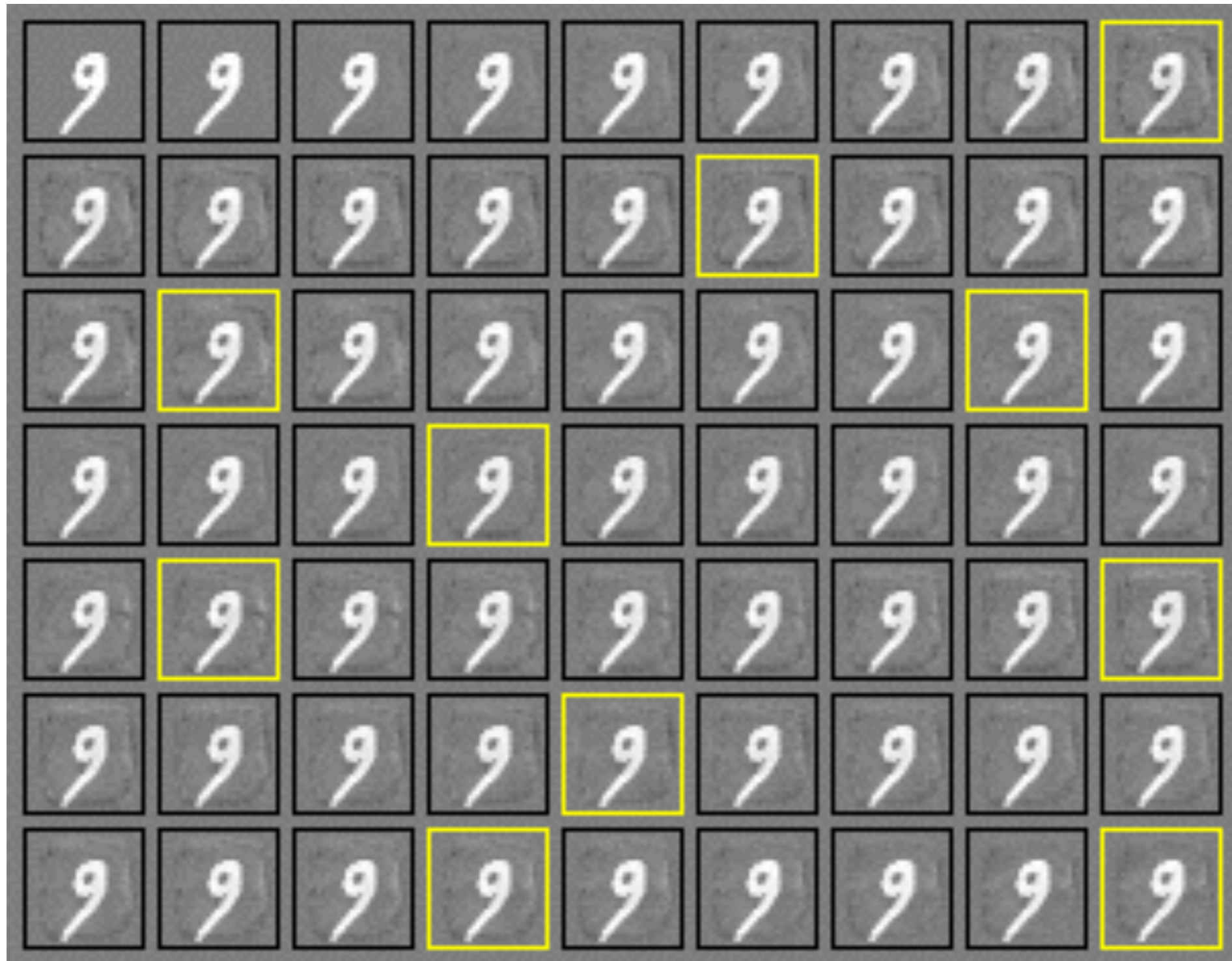
Szegedy et al 2013: fool ImageNet classifiers imperceptibly

Goodfellow et al 2014: cheap, closed form attack

Turning Objects into “Airplanes”



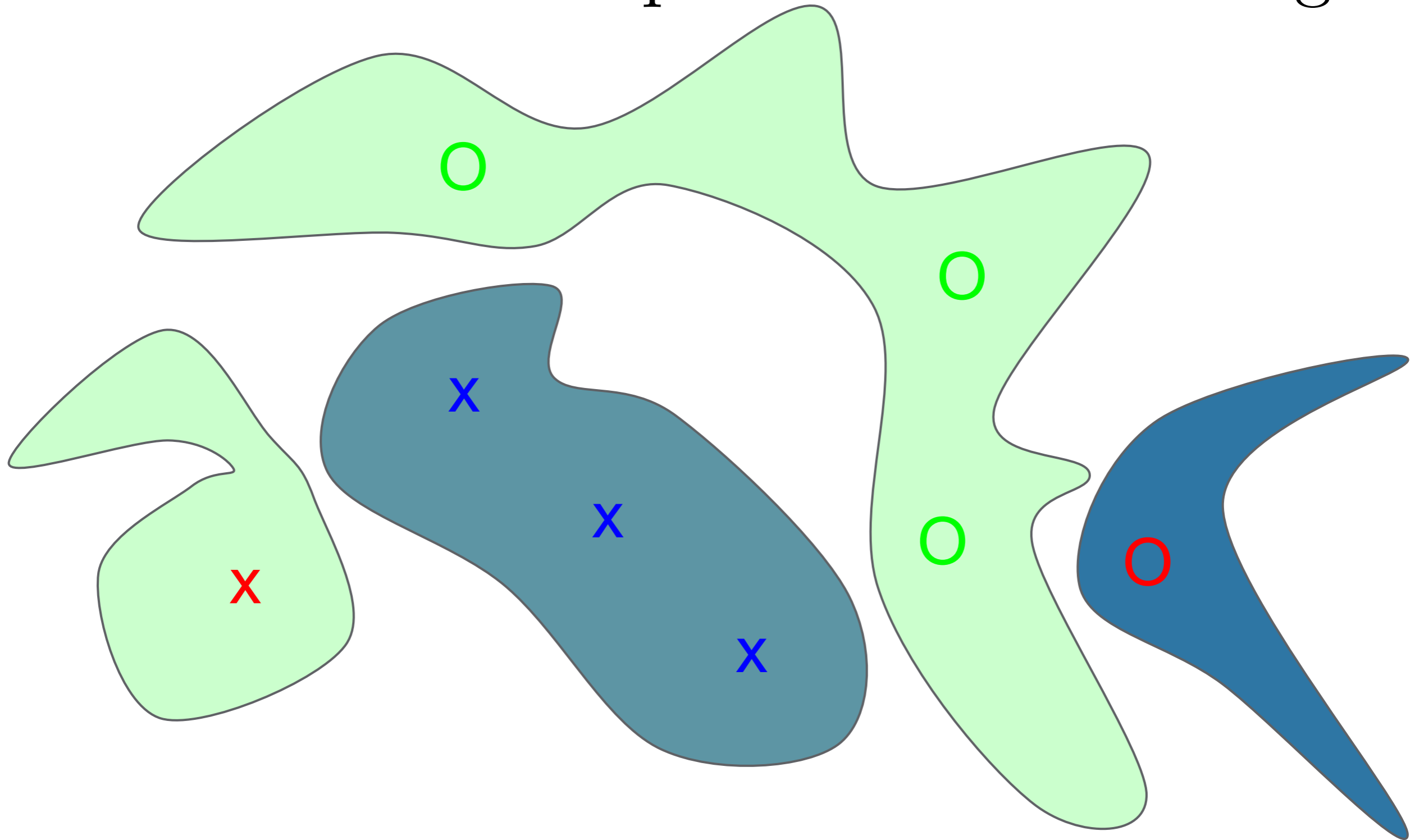
Attacking a Linear Model



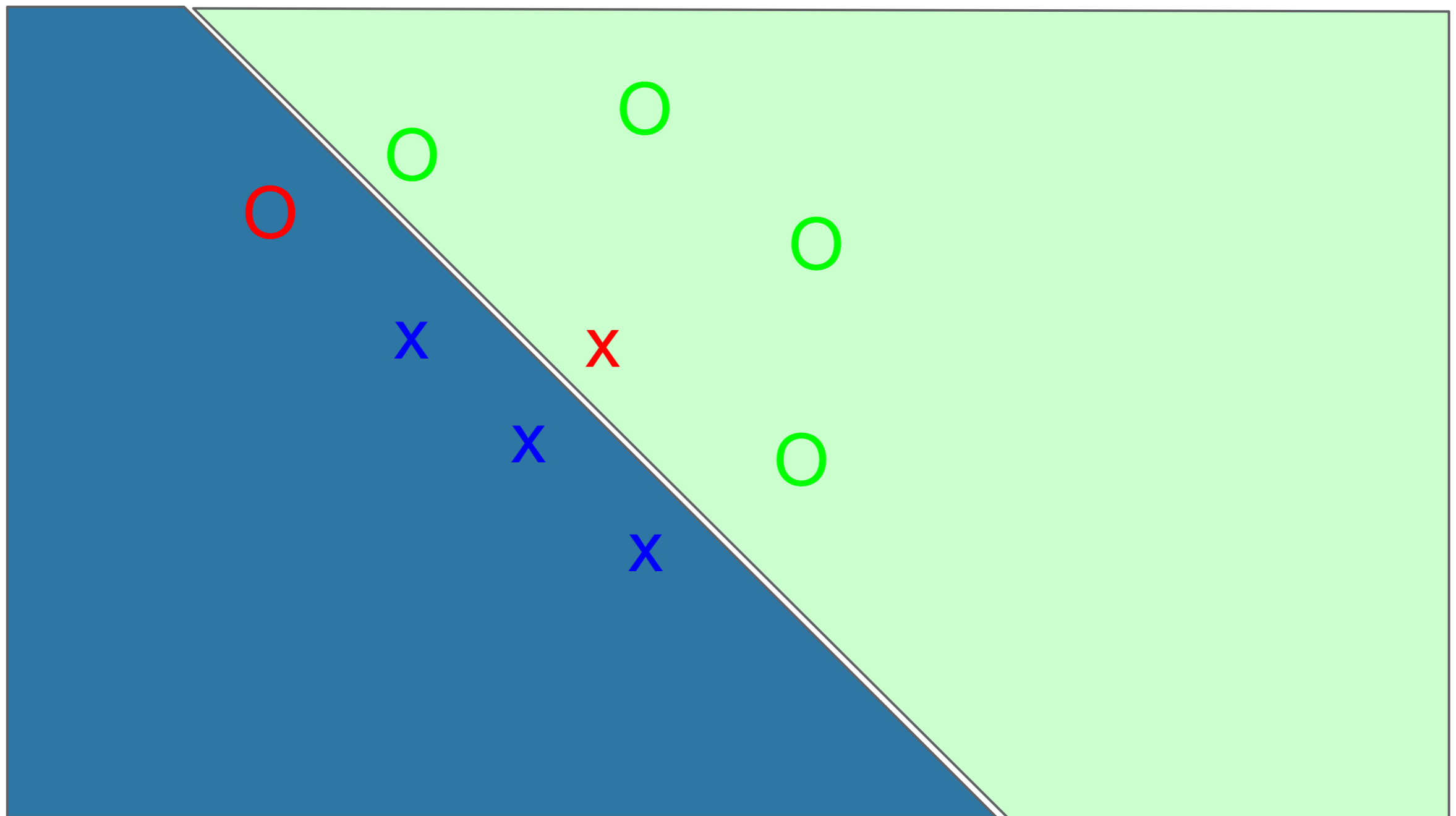
Not just for neural nets

- Linear models
 - Logistic regression
 - Softmax regression
 - SVMs
- Decision trees
- Nearest neighbors

Adversarial Examples from Overfitting

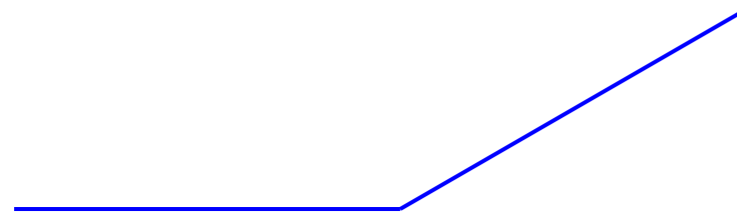


Adversarial Examples from Excessive Linearity

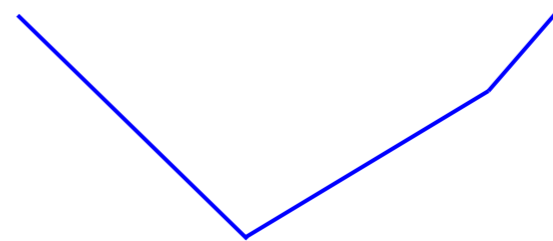


Modern deep nets are very piecewise linear

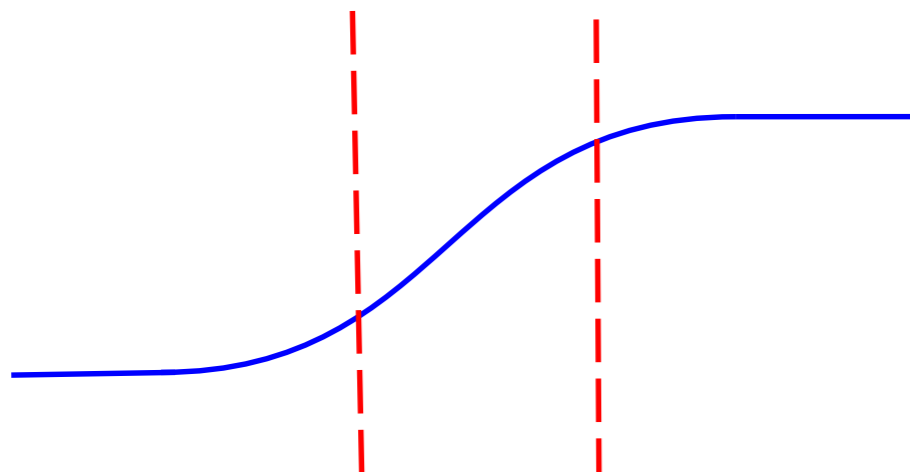
Rectified linear unit



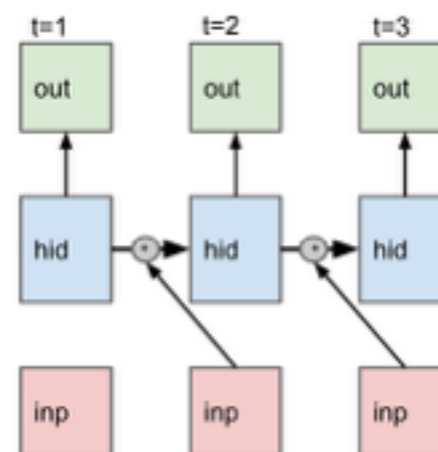
Maxout



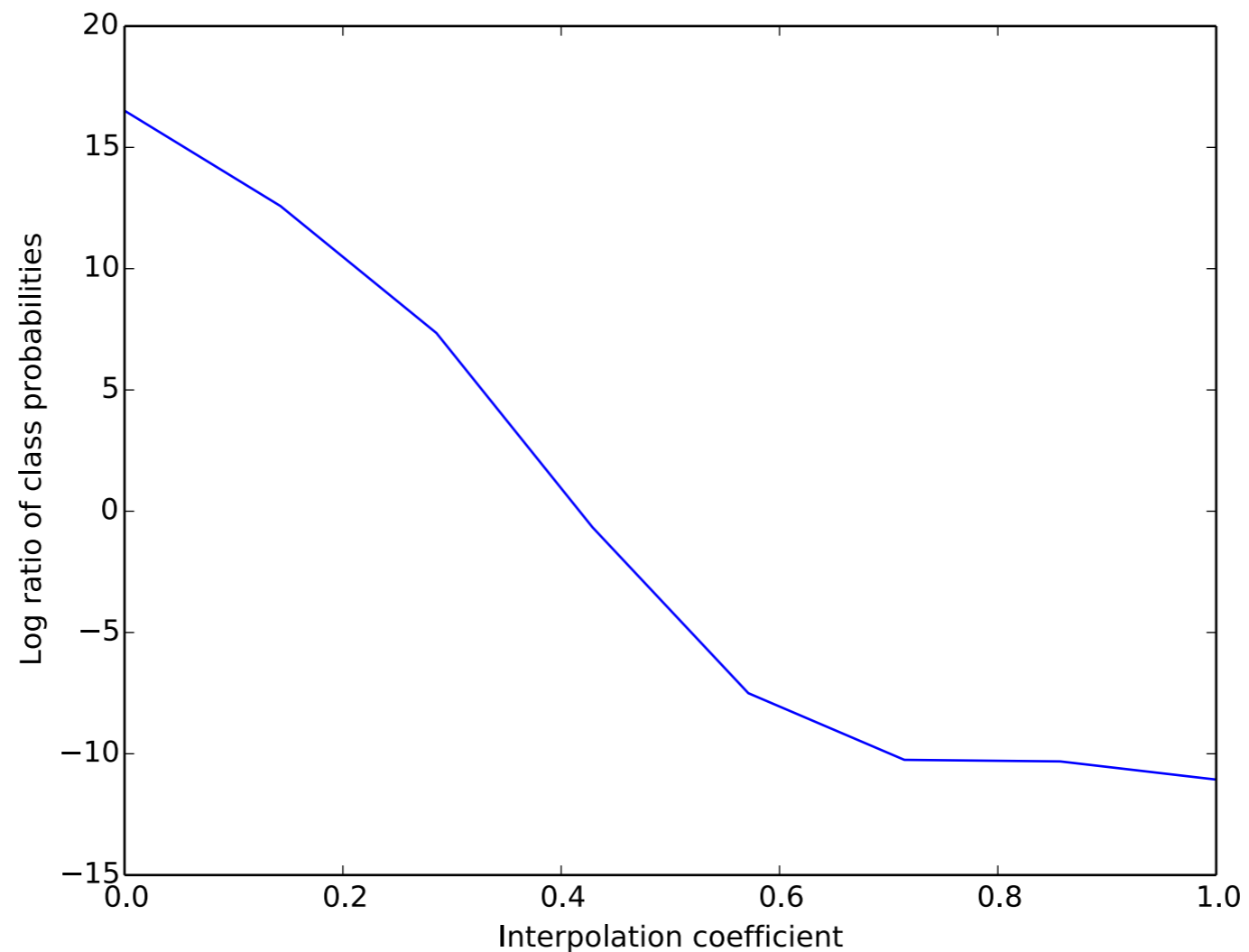
Carefully tuned sigmoid



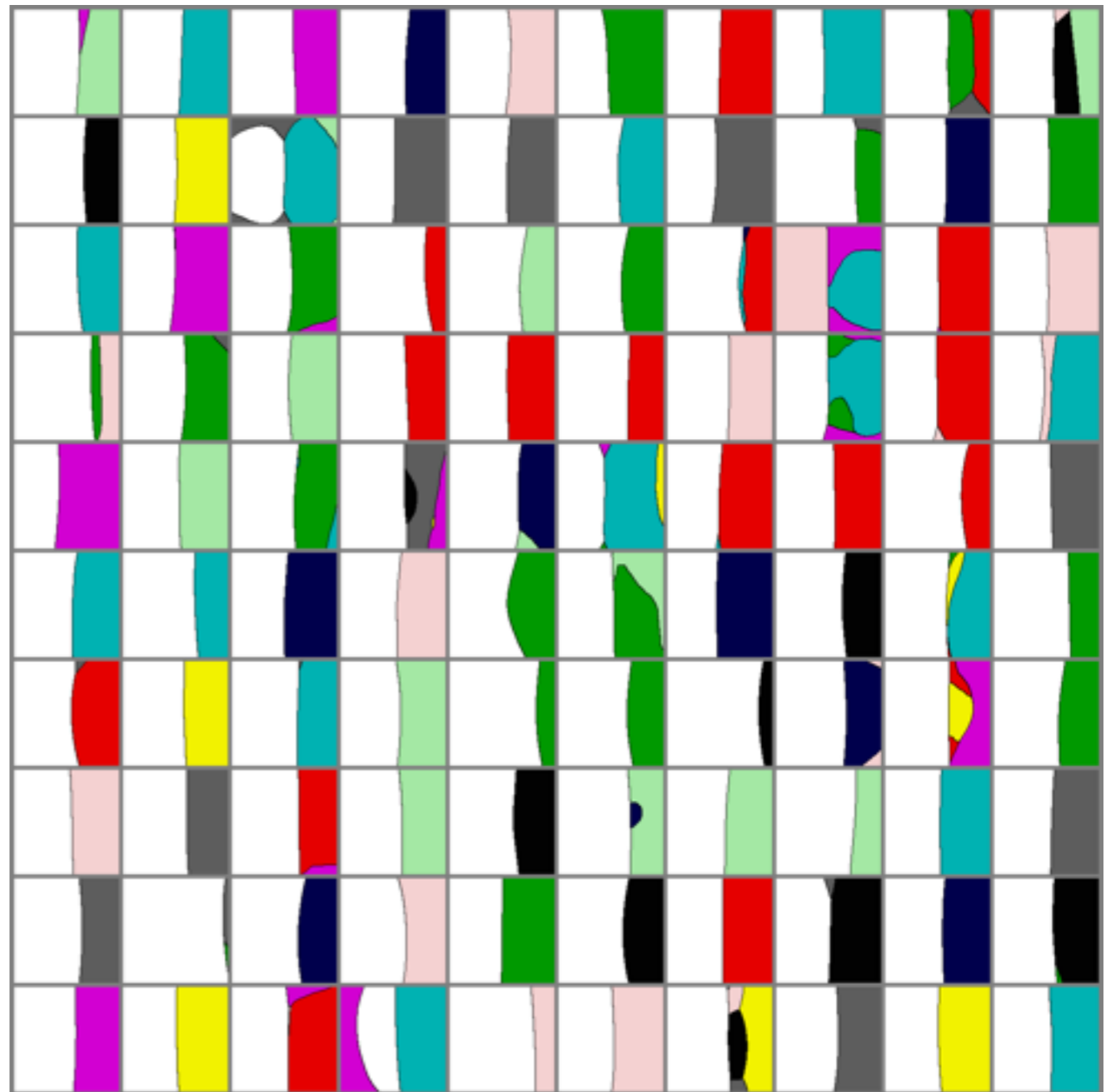
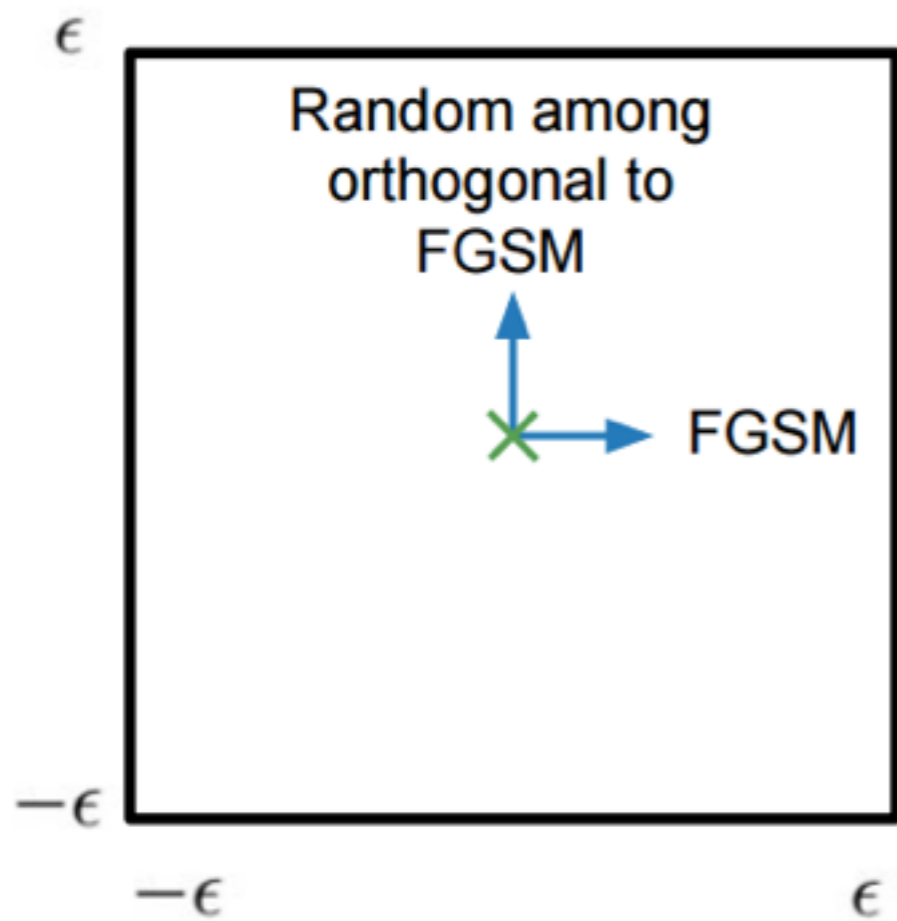
LSTM



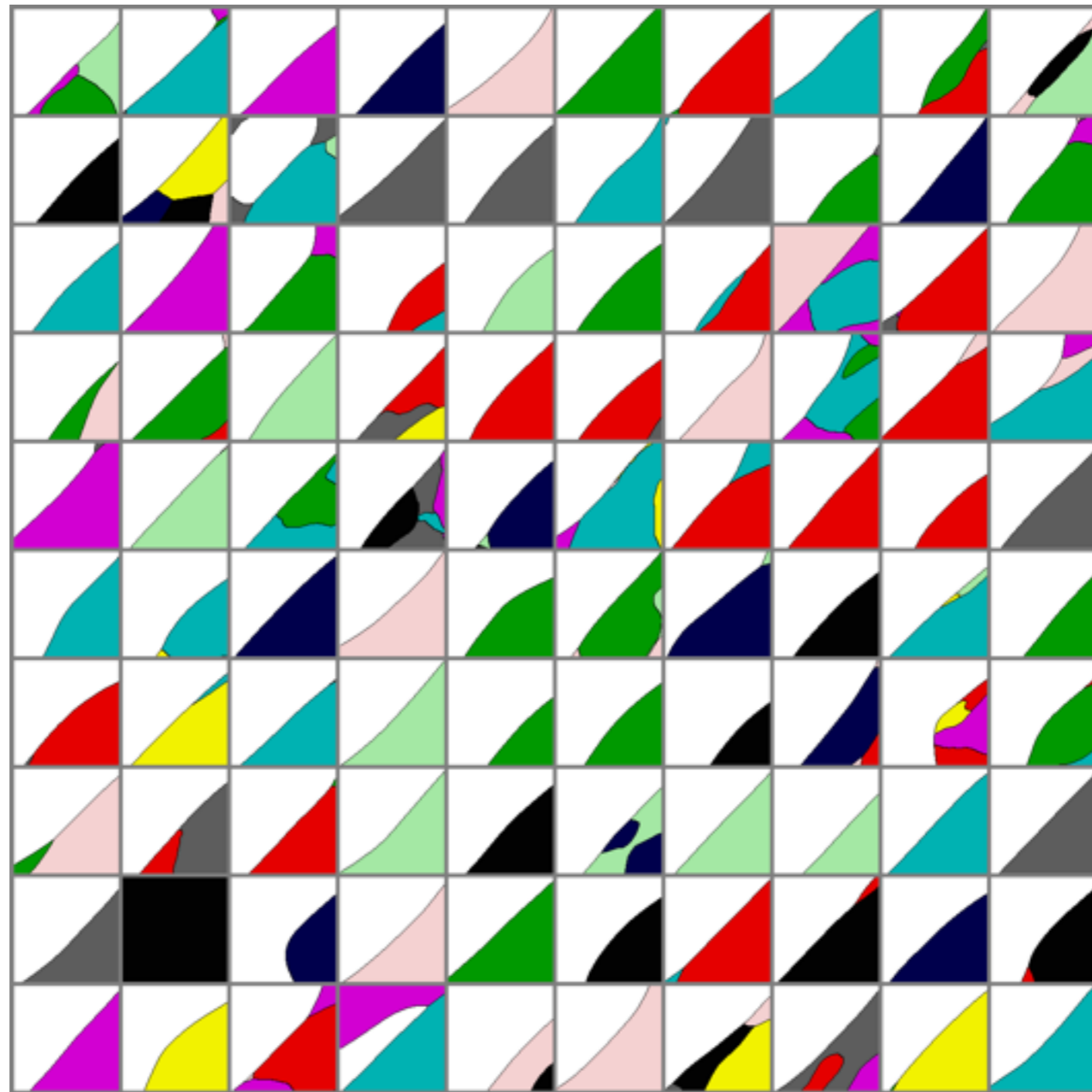
Nearly Linear Responses in Practice



Maps of Adversarial and Random Cross-Sections

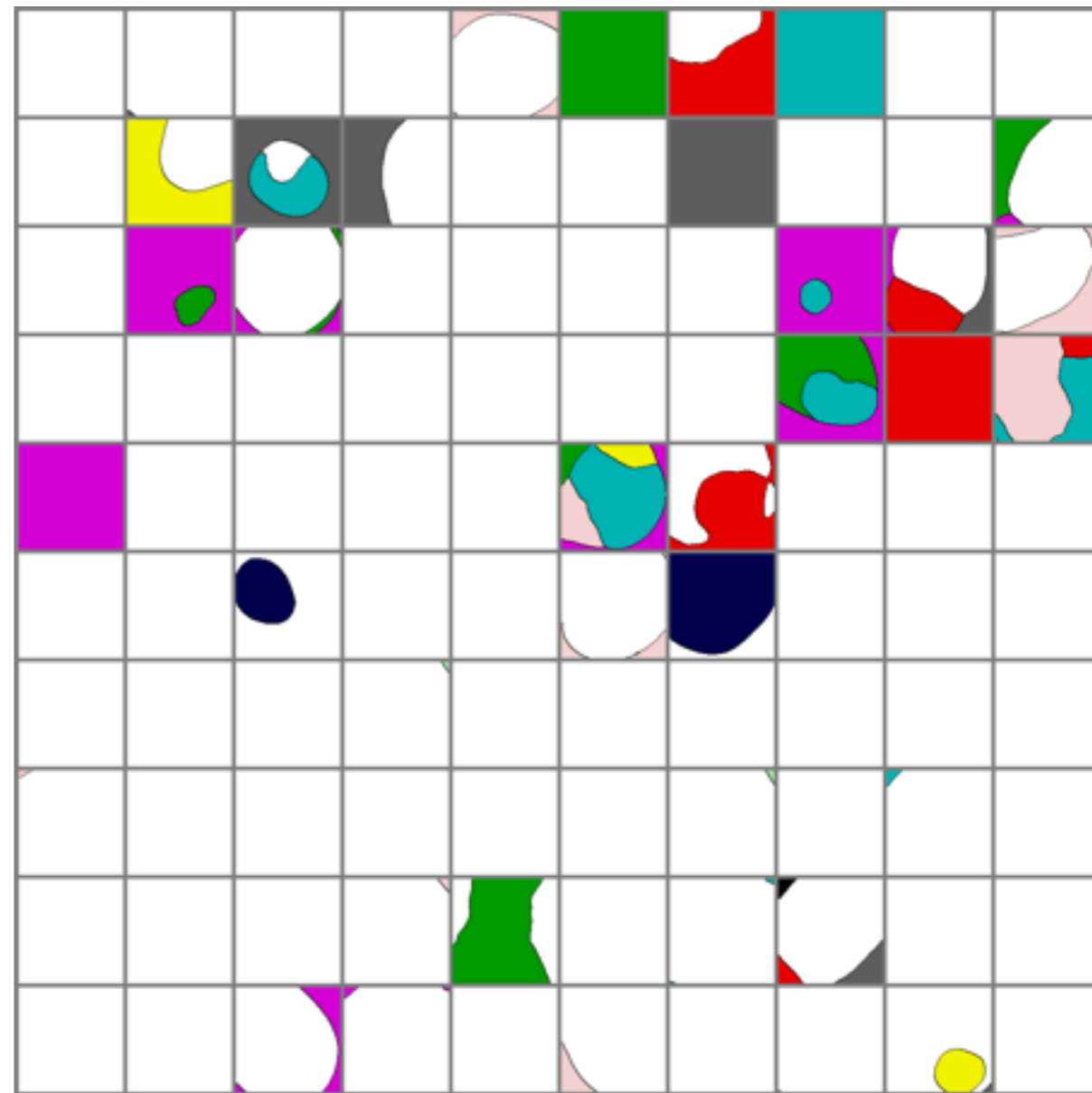
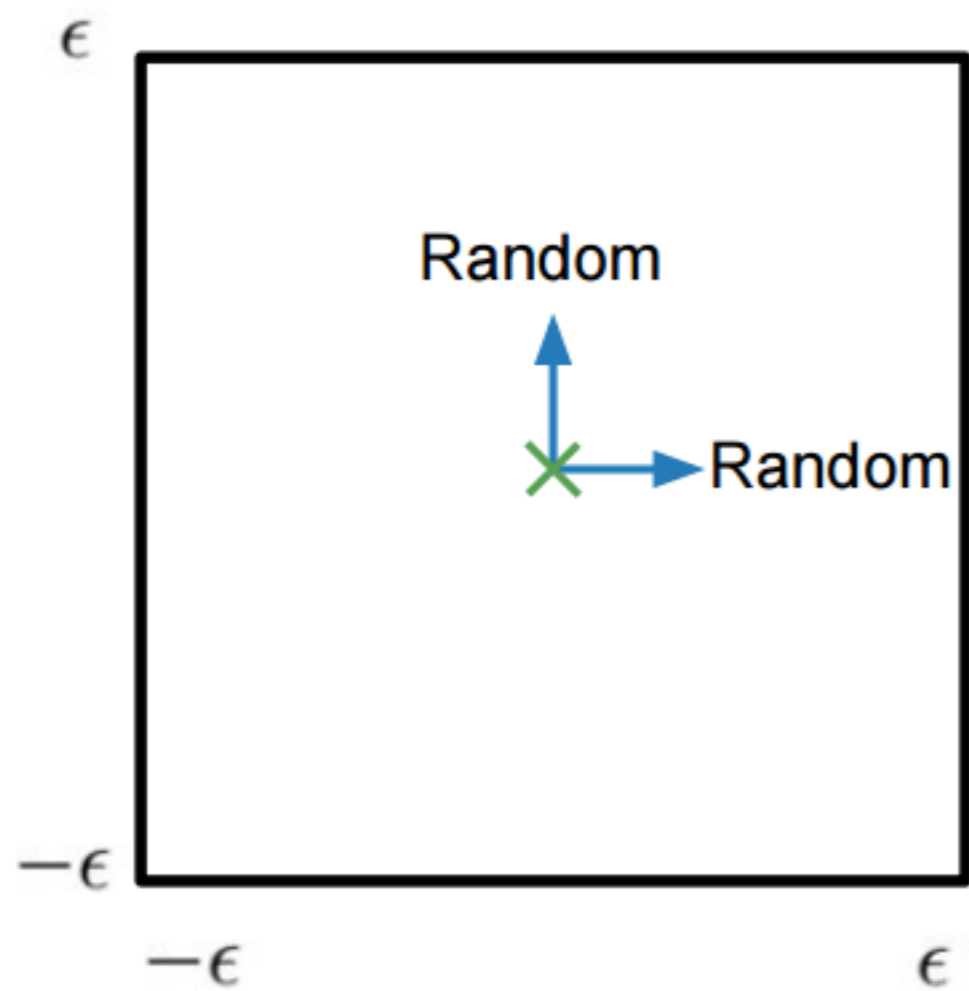


Maps of Adversarial Cross-Sections



Maps of Random Cross-Sections

Adversarial examples
are not noise



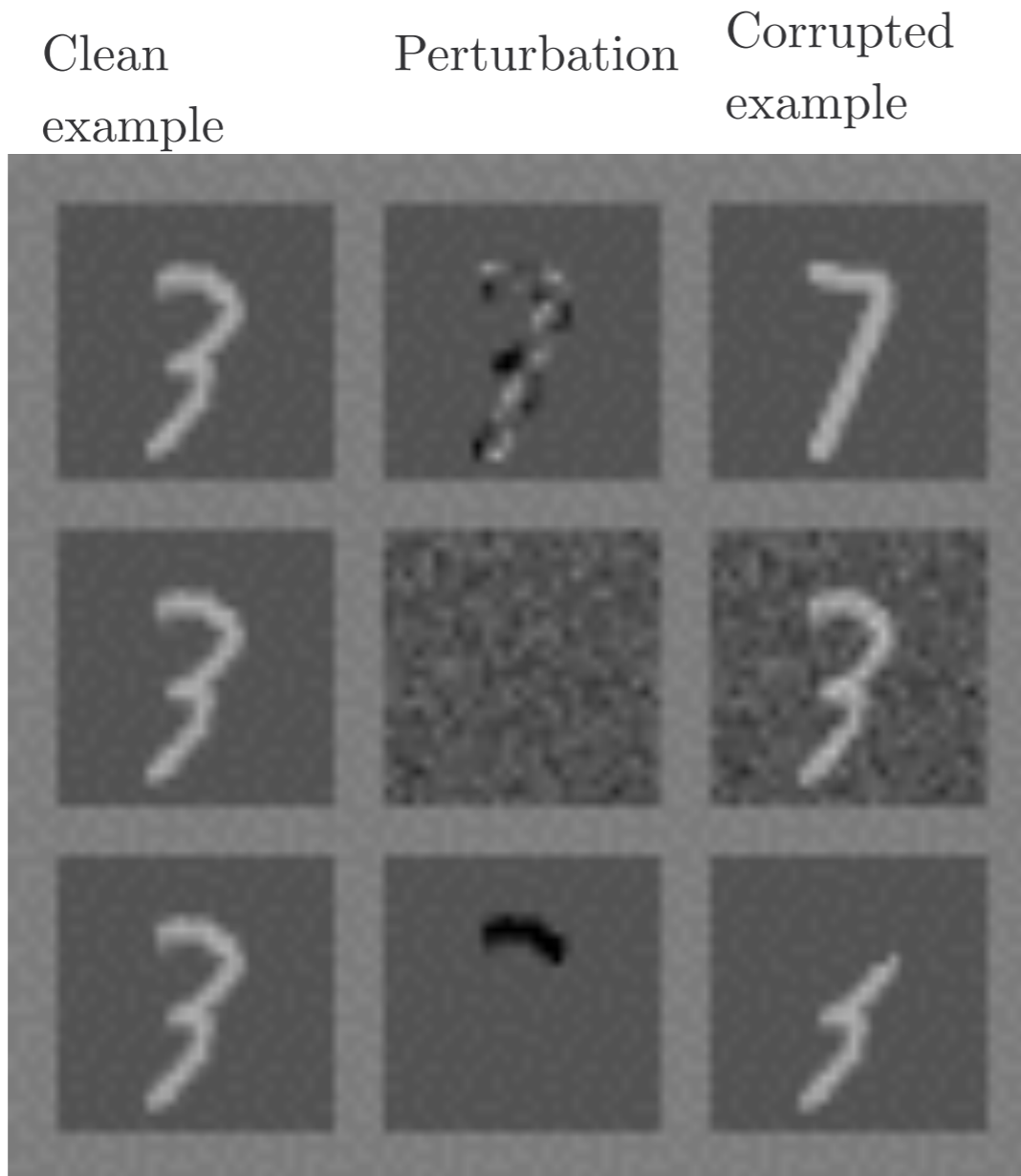
Clever Hans



“Clever Hans,
Clever
Algorithms,”
Bob Sturm)



Small inter-class distances



Perturbation changes the true class

Random perturbation does not change the class

Perturbation changes the input to “rubbish class”

All three perturbations have L2 norm 3.96

This is actually small. We typically use 7!

The Fast Gradient Sign Method

$$J(\tilde{\mathbf{x}}, \boldsymbol{\theta}) \approx J(\mathbf{x}, \boldsymbol{\theta}) + (\tilde{\mathbf{x}} - \mathbf{x})^\top \nabla_{\mathbf{x}} J(\mathbf{x}).$$

Maximize

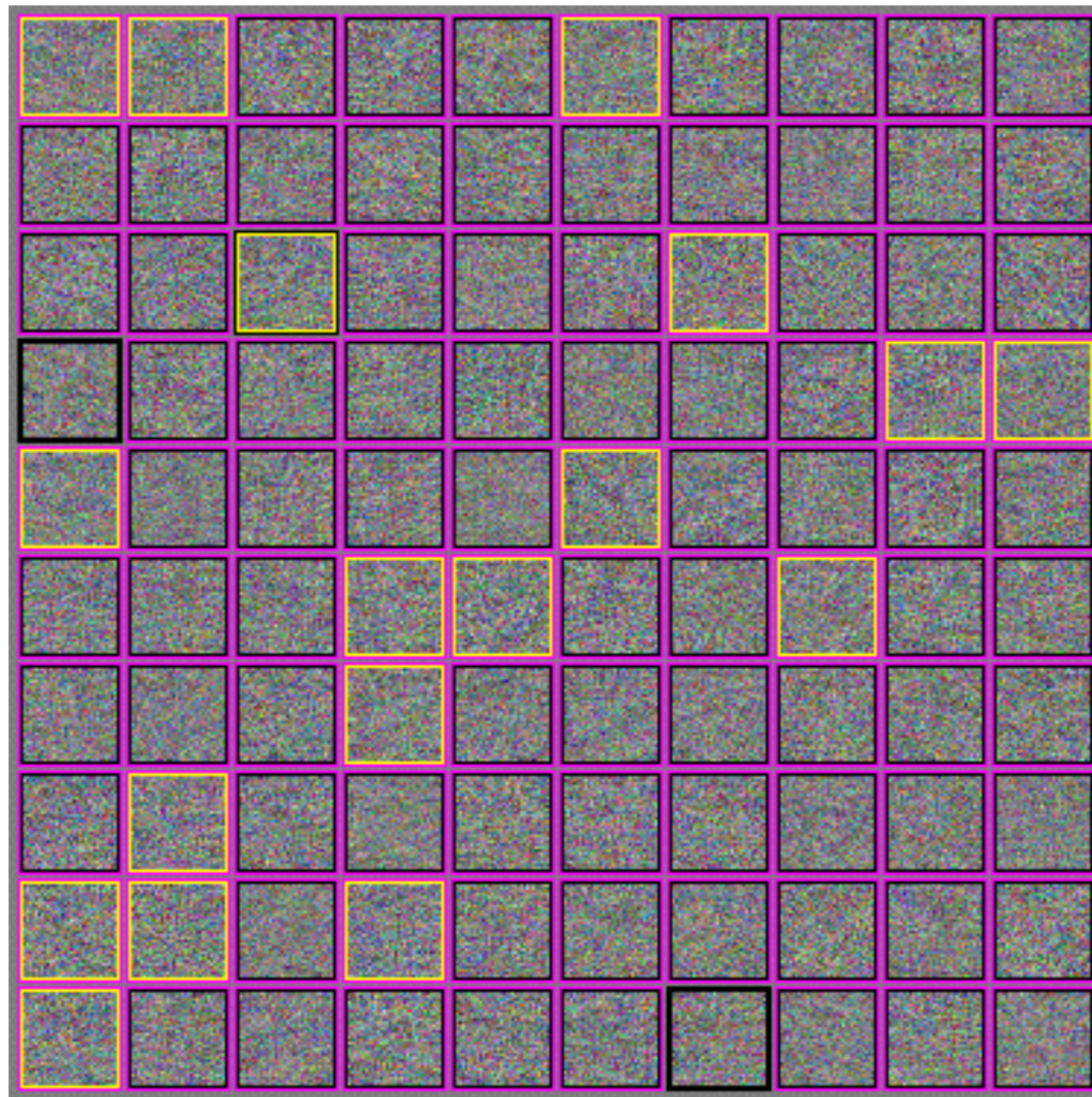
$$J(\mathbf{x}, \boldsymbol{\theta}) + (\tilde{\mathbf{x}} - \mathbf{x})^\top \nabla_{\mathbf{x}} J(\mathbf{x})$$

subject to

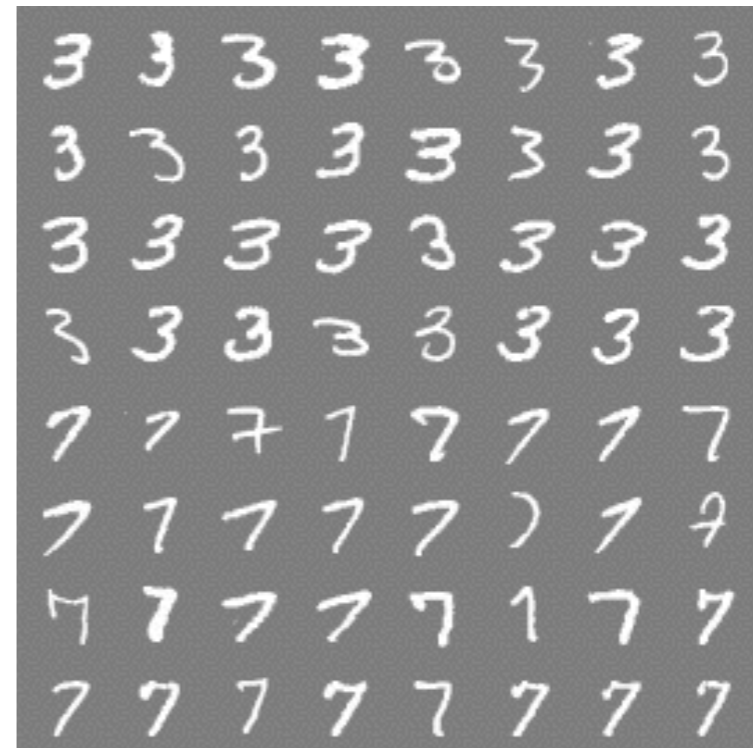
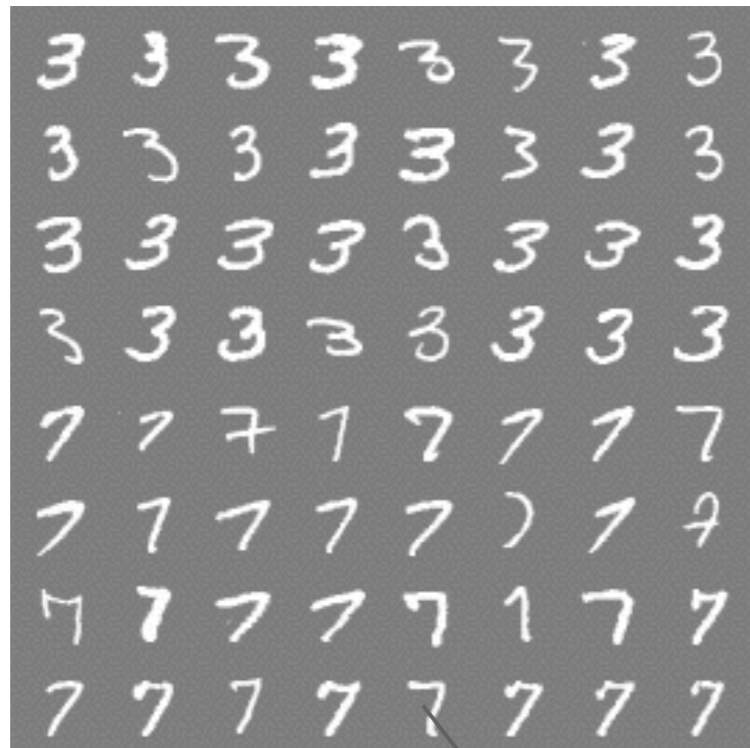
$$\|\tilde{\mathbf{x}} - \mathbf{x}\|_\infty \leq \epsilon$$

$$\Rightarrow \tilde{\mathbf{x}} = \mathbf{x} + \epsilon \text{sign}(\nabla_{\mathbf{x}} J(\mathbf{x})).$$

Wrong almost everywhere



Cross-model, cross-dataset generalization



Cross-technique transferability

Source Machine Learning Technique	DNN	LR	SVM	DT	kNN	Ens.
DNN	38.27	23.02	64.32	79.31	8.36	20.72
LR	6.31	91.64	91.43	87.42	11.29	44.14
SVM	2.51	36.56	100.0	80.03	5.19	15.67
DT	0.82	12.22	8.85	89.29	3.31	5.11
kNN	11.75	42.89	82.16	82.95	41.65	31.92

(Papernot 2016)

Transferability Attack

Target model with unknown weights, machine learning algorithm, training set; maybe non-differentiable

Train your own model

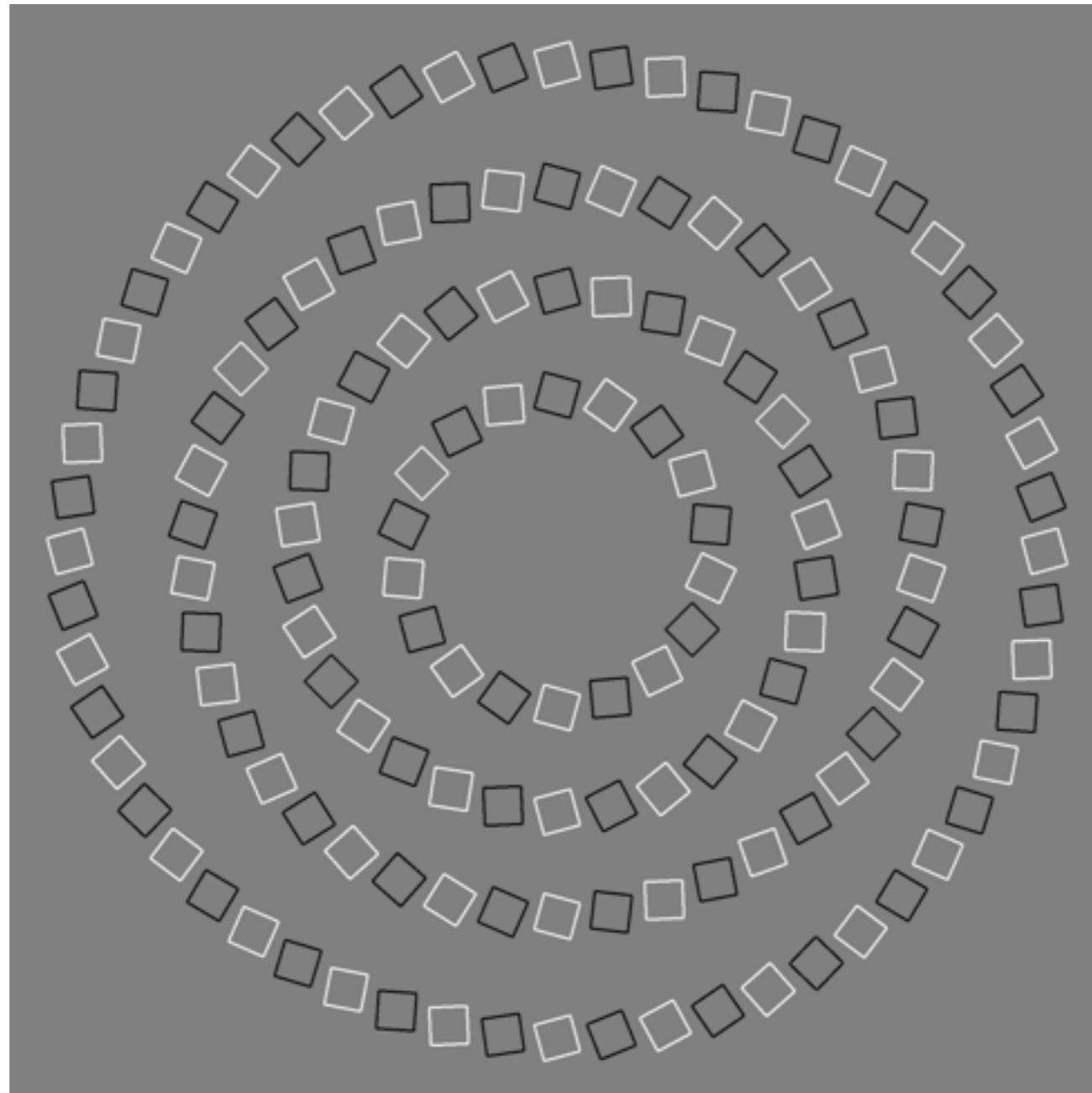
Substitute model mimicking target model with known, differentiable function

Adversarial crafting against substitute

Adversarial examples

Deploy adversarial examples against the target; transferability property results in them succeeding

Adversarial Examples in the Human Brain



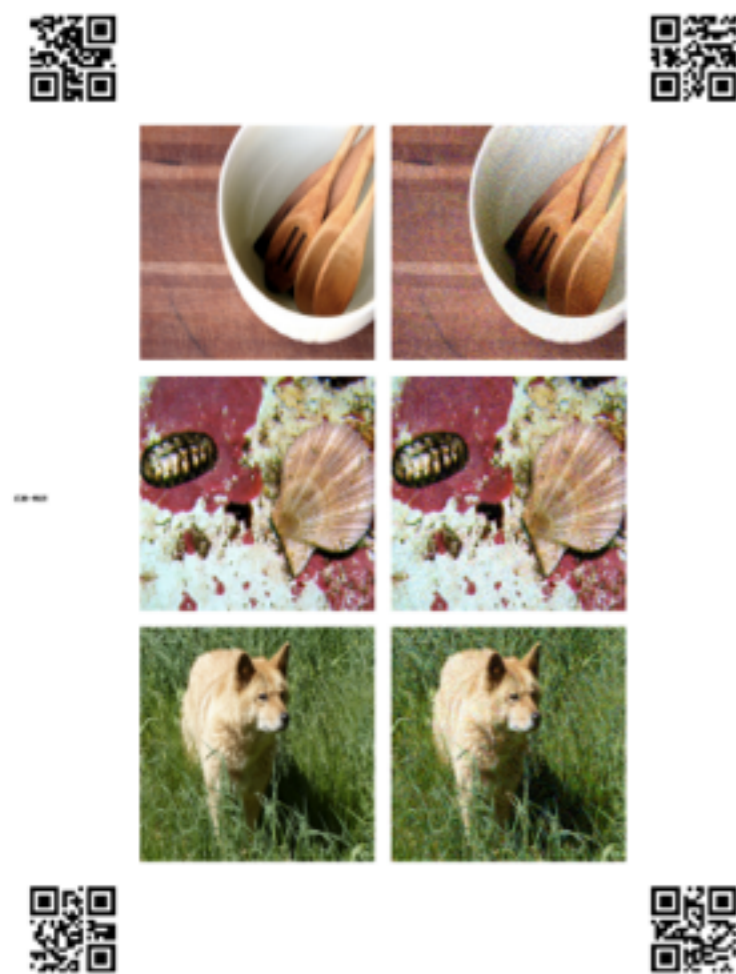
These are
concentric
circles,
not
intertwined
spirals.

(Pinna and Gregory, 2002)

Practical Attacks

- Fool real classifiers trained by remotely hosted API (MetaMind, Amazon, Google)
- Fool malware detector networks
- Display adversarial examples in the physical world and fool machine learning systems that perceive them through a camera

Adversarial Examples in the Physical World



(a) Printout



(b) Photo of printout



(c) Cropped image

Failed defenses

Generative
pretraining

Removing perturbation
with an autoencoder

Adding noise
at test time

Ensembles

Confidence-reducing
perturbation at test time

Error correcting
codes

Multiple glimpses

Weight decay

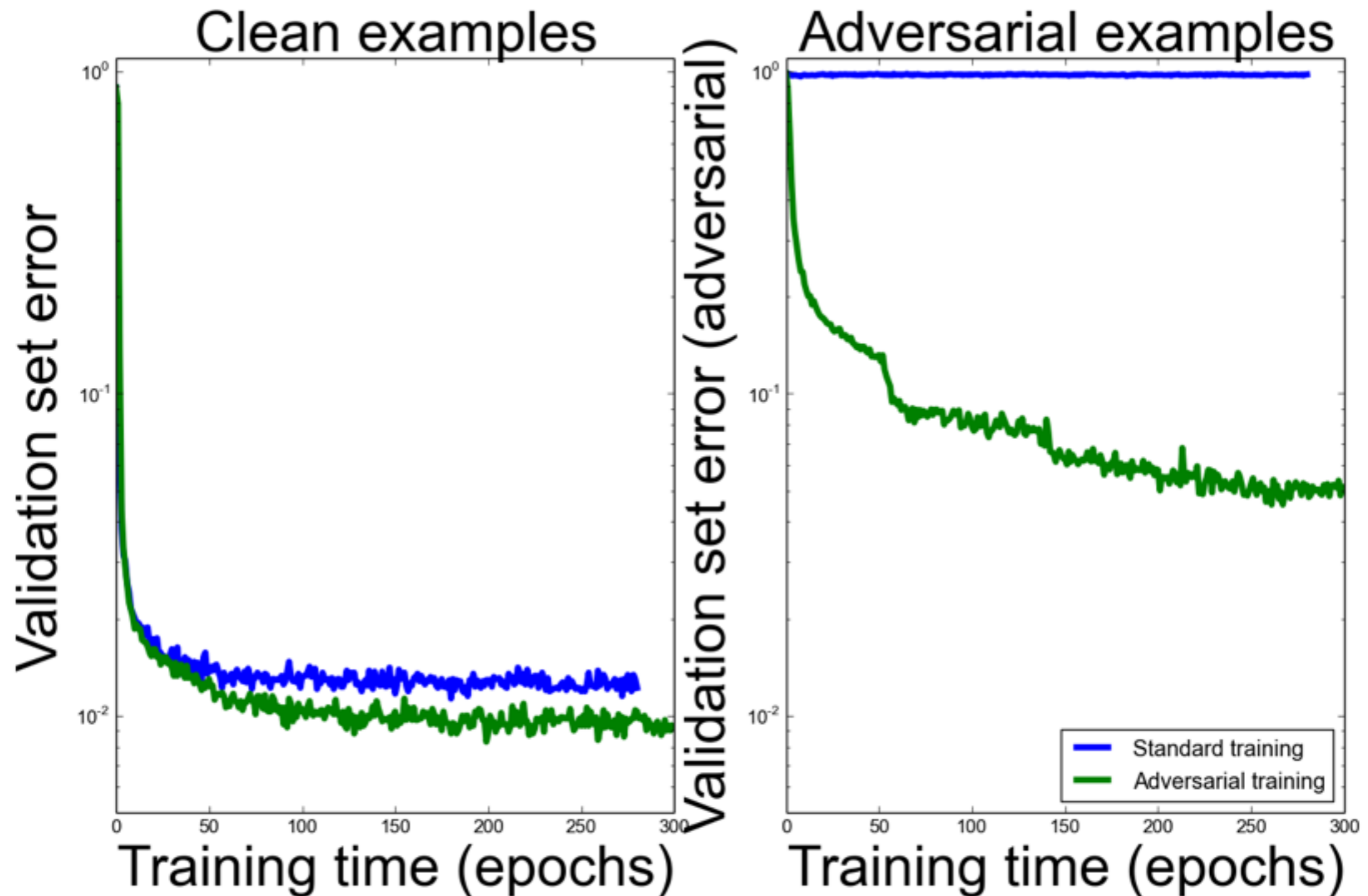
Double backprop

Adding noise
at train time

Various
non-linear units

Dropout

Training on Adversarial Examples



Adversarial Training

Labeled as bird



Still has same label (bird)



Decrease
probability
of bird class

Virtual Adversarial Training

Unlabeled; model
guesses it's probably
a bird, maybe a plane



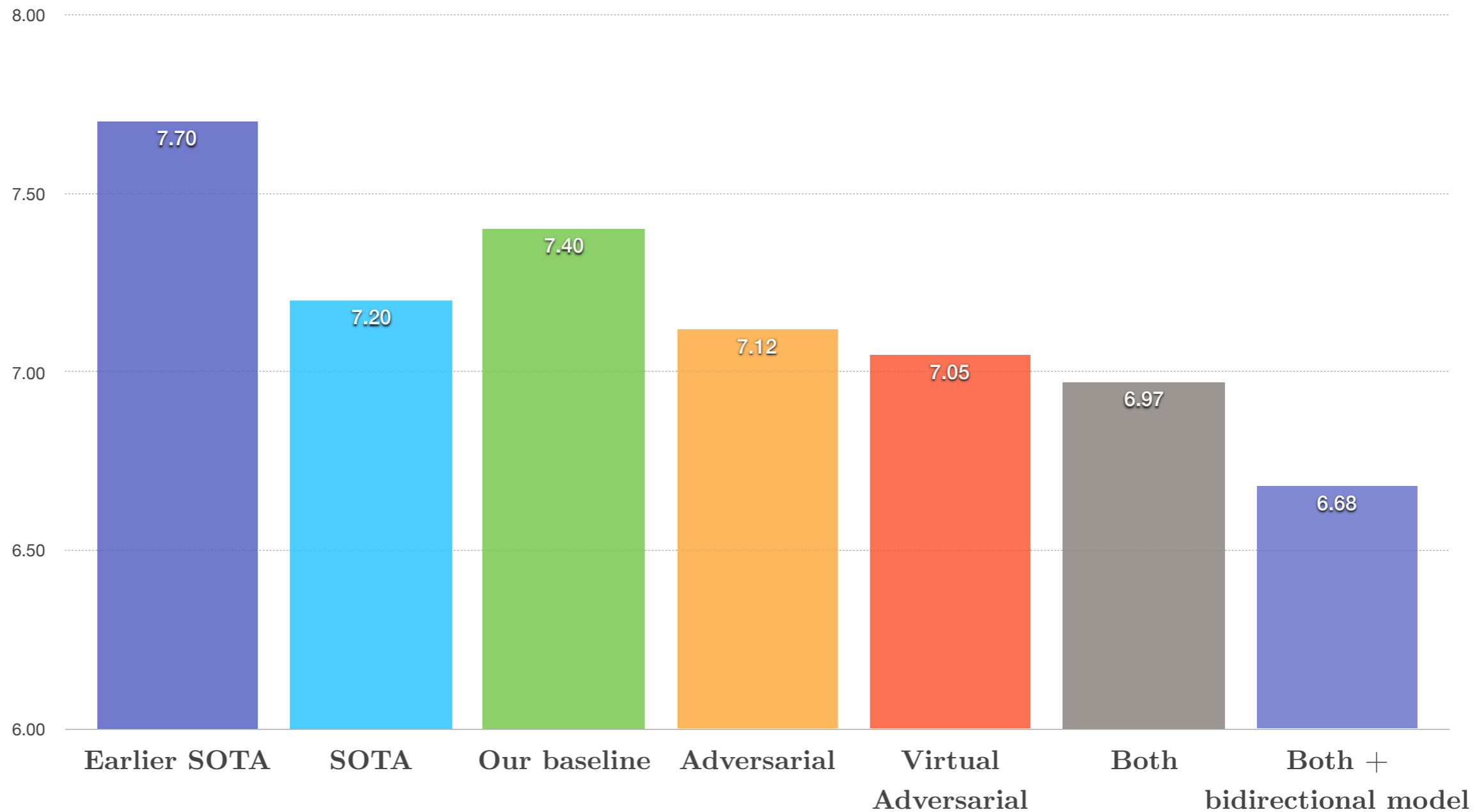
New guess should
match old guess
(probably bird, maybe plane)



→
Adversarial
perturbation
intended to
change the guess

Text Classification with VAT

RCV1 Misclassification Rate



Zoomed in for legibility

Conclusion

- Attacking is easy
- Defending is difficult
- Benchmarking vulnerability is training
- Adversarial training provides regularization and semi-supervised learning