

Thermometer Encoding: One Hot Way to Resist Adversarial Examples

Stanford, 2017-11-16



Jacob
Buckman*



Aurko Roy*

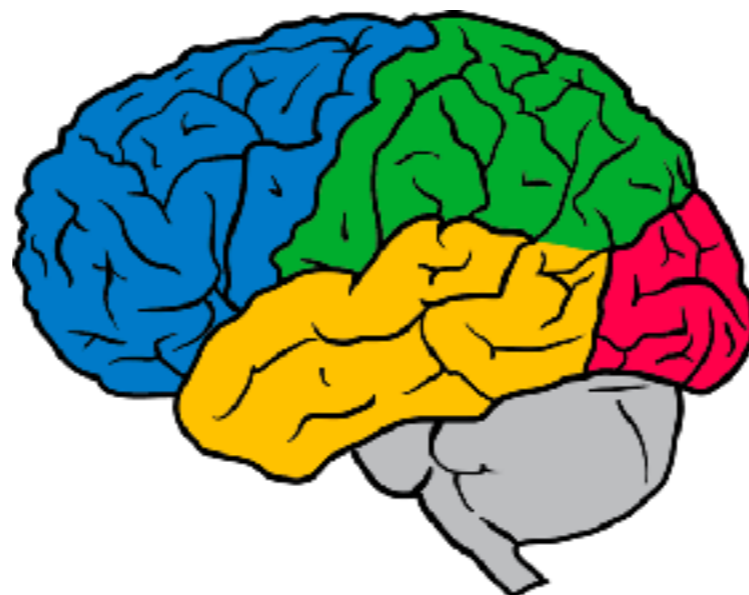


Colin Raffel



Ian
Goodfellow

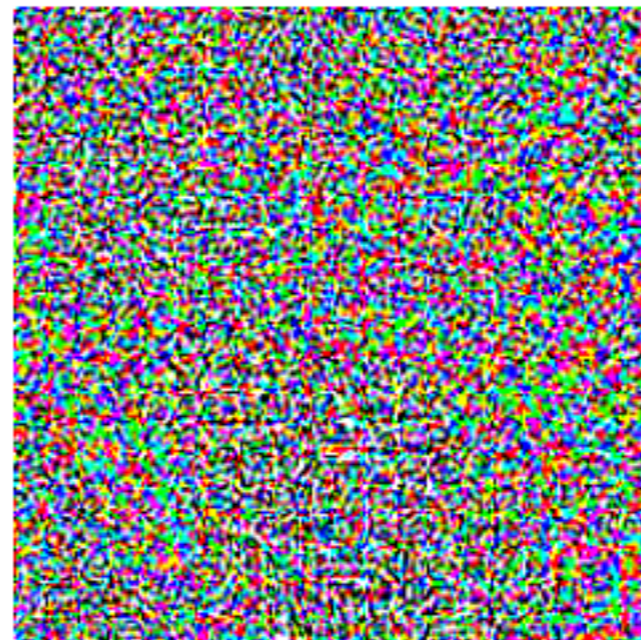
*joint first author



Adversarial Examples



+ .007 ×



=

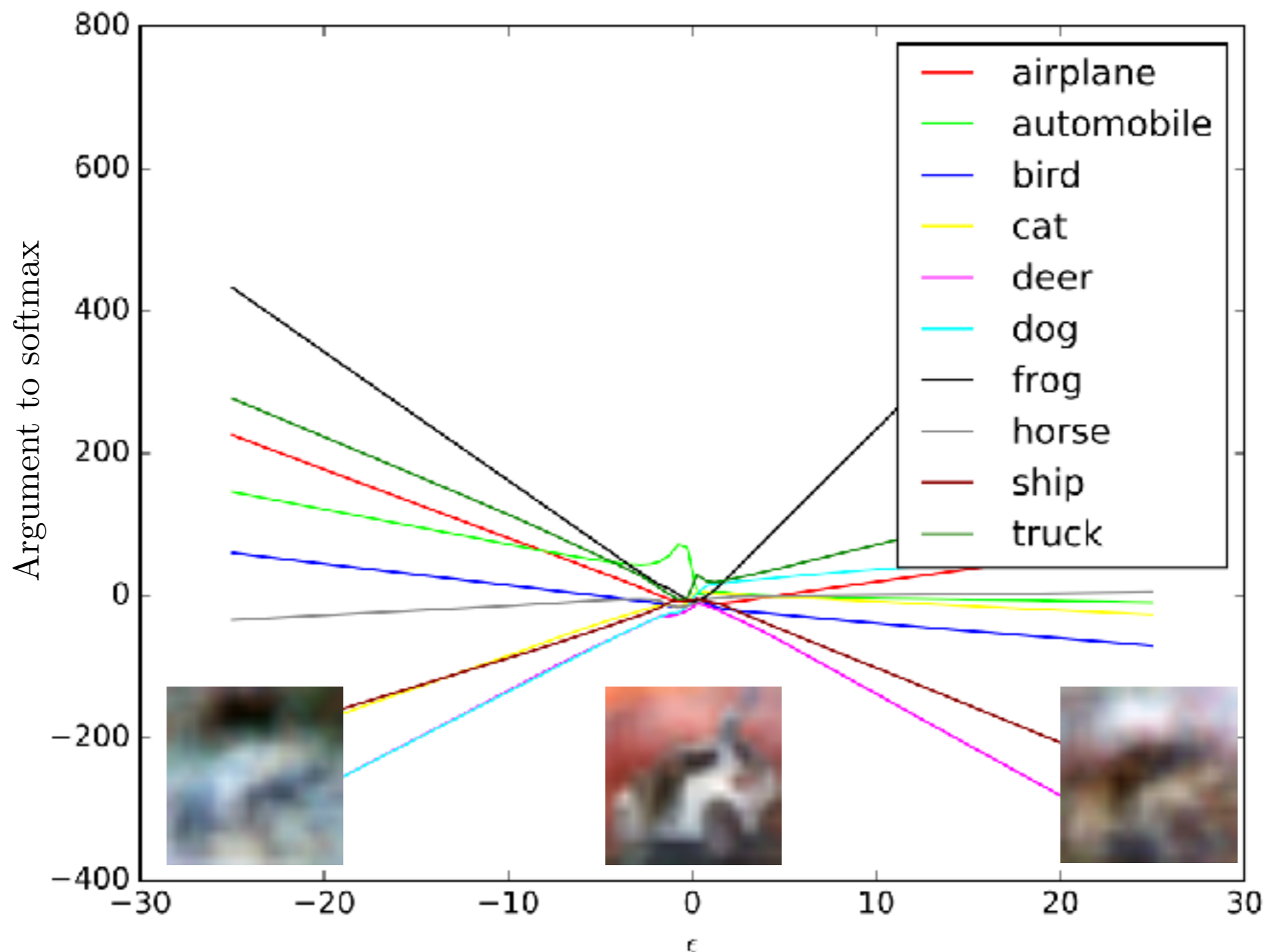


Probably panda

Adversarial
perturbation

Definitely
gibbon

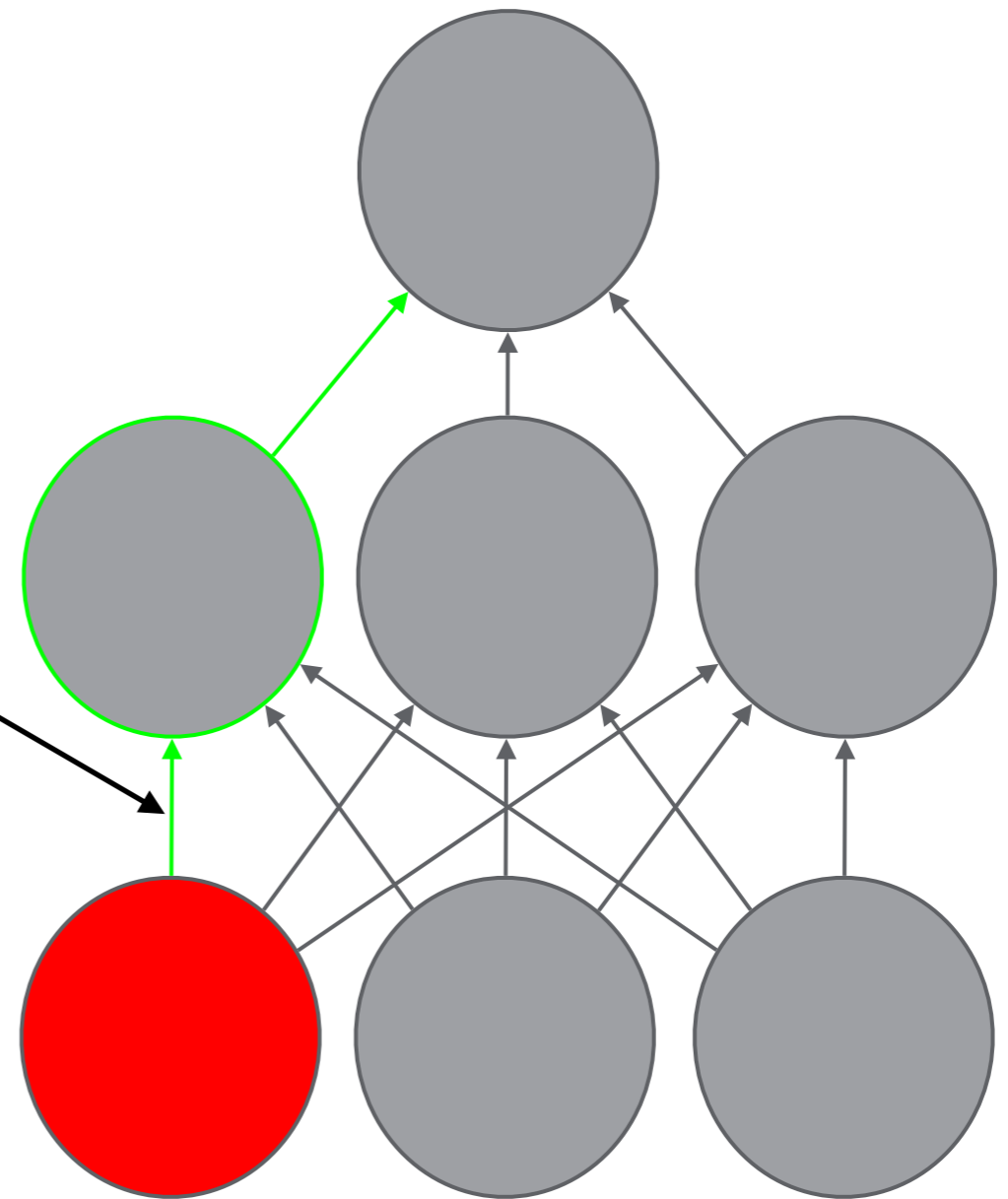
Unreasonable Linear Extrapolation



Difficult to train extremely nonlinear hidden layers

To train:

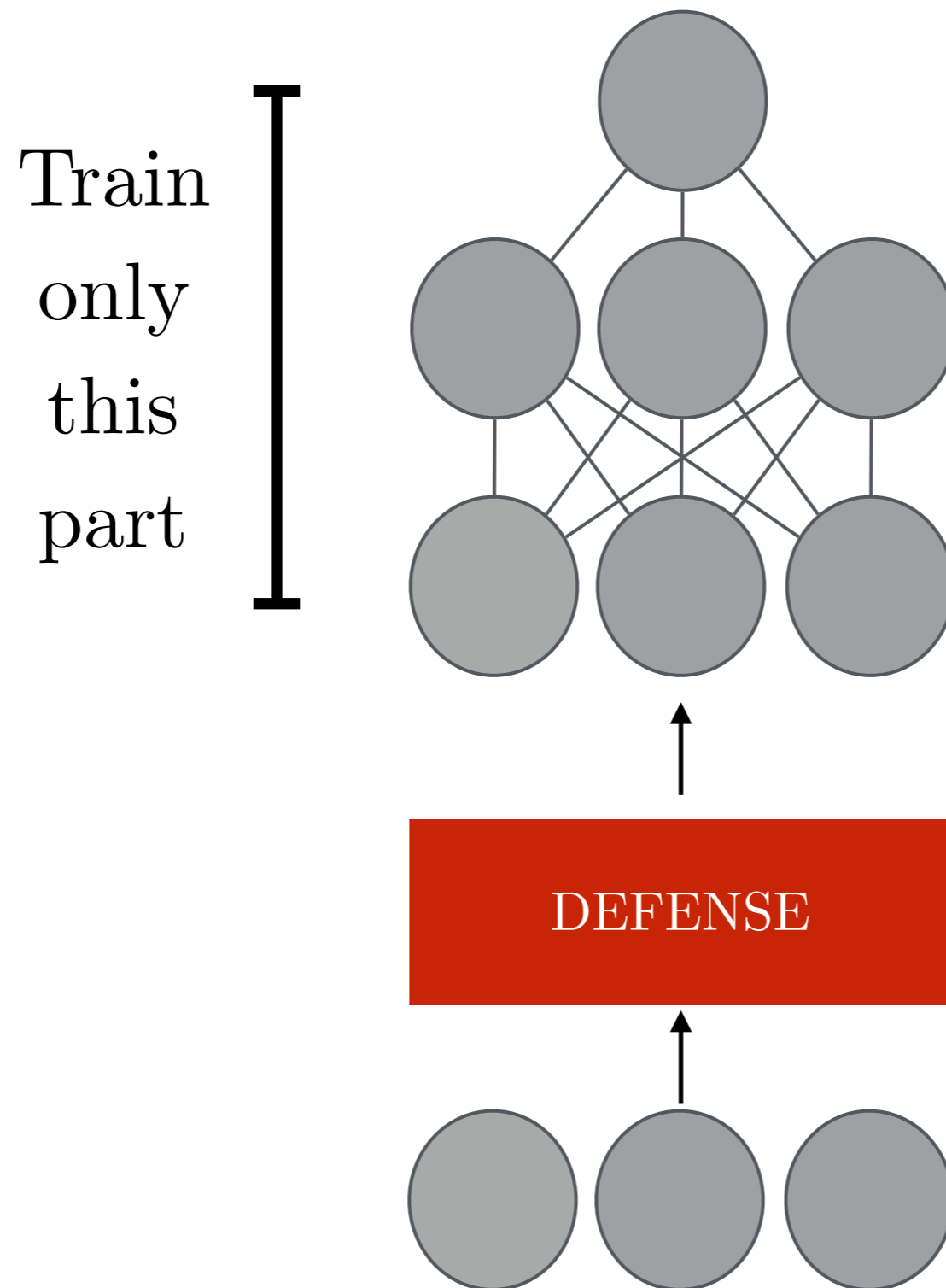
changing this weight needs to
have a large, predictable effect



To defend:

changing this input needs
to have a small or
unpredictable effect

Idea: edit only the input layer



Real-valued

Quantized

0.13

0.15

0.66

0.65

0.92

0.95

Discretized (one-hot)

Discretized (thermometer)

[0100000000]

[0111111111]

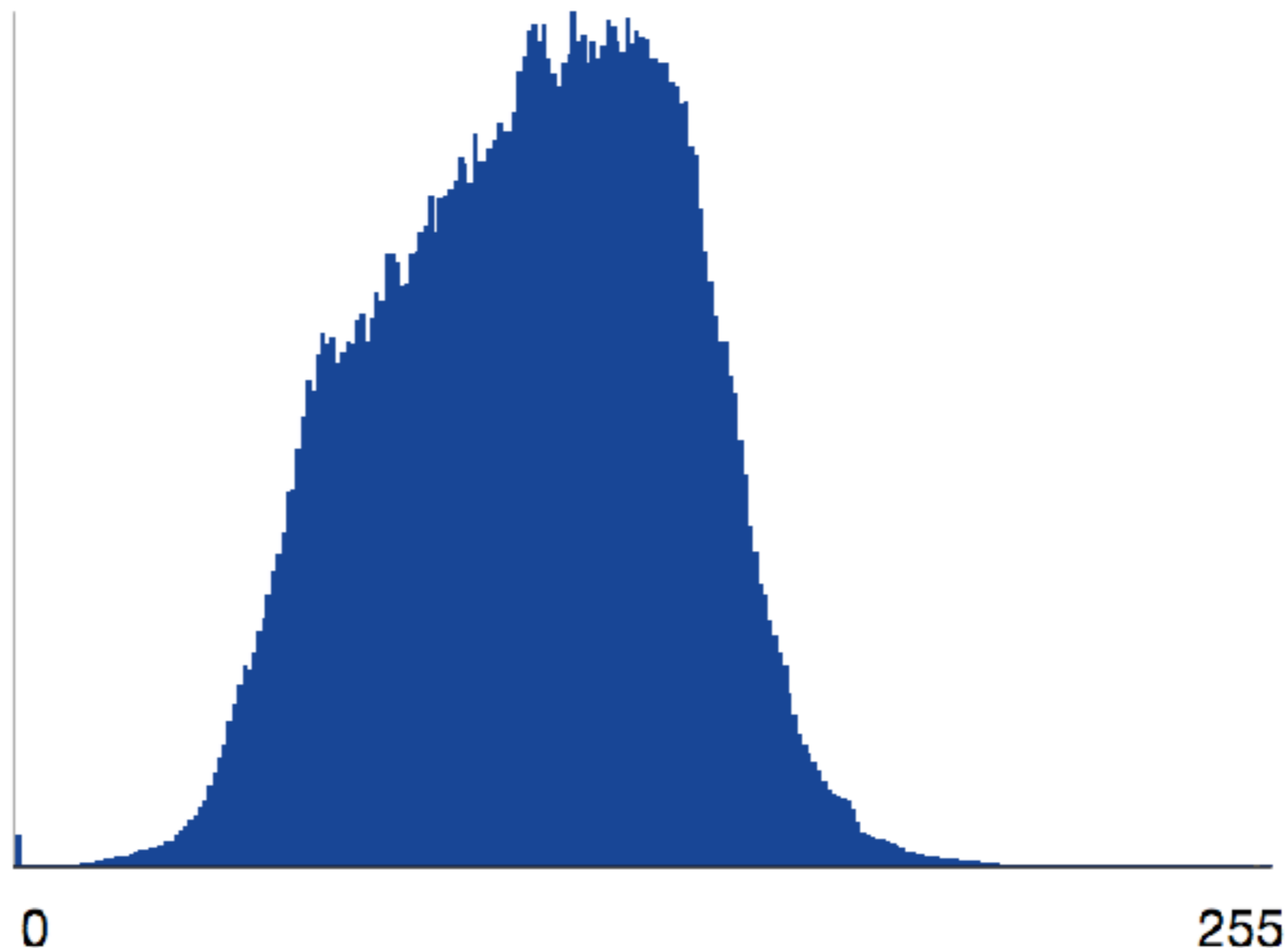
[0000001000]

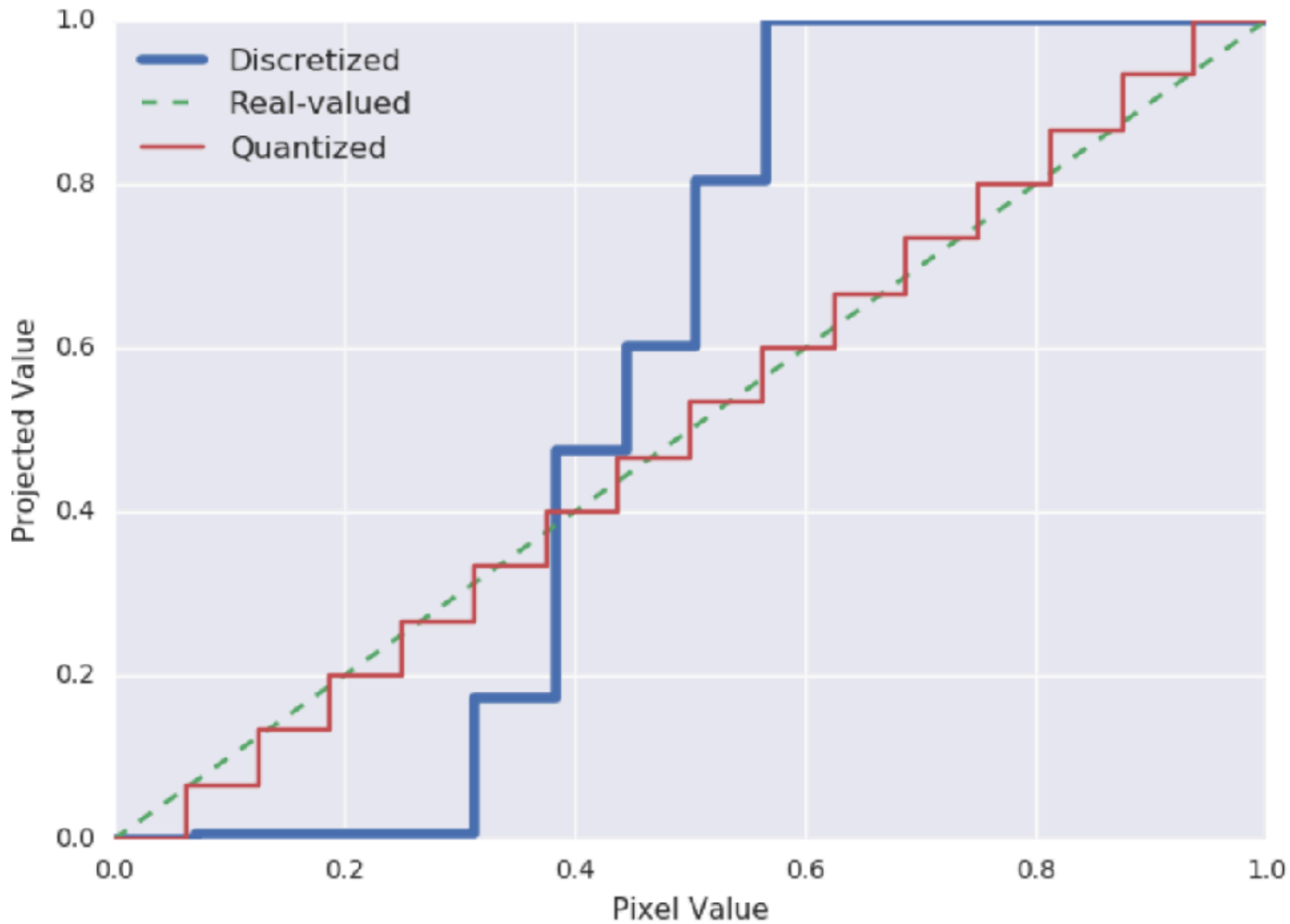
[0000001111]

[0000000001]

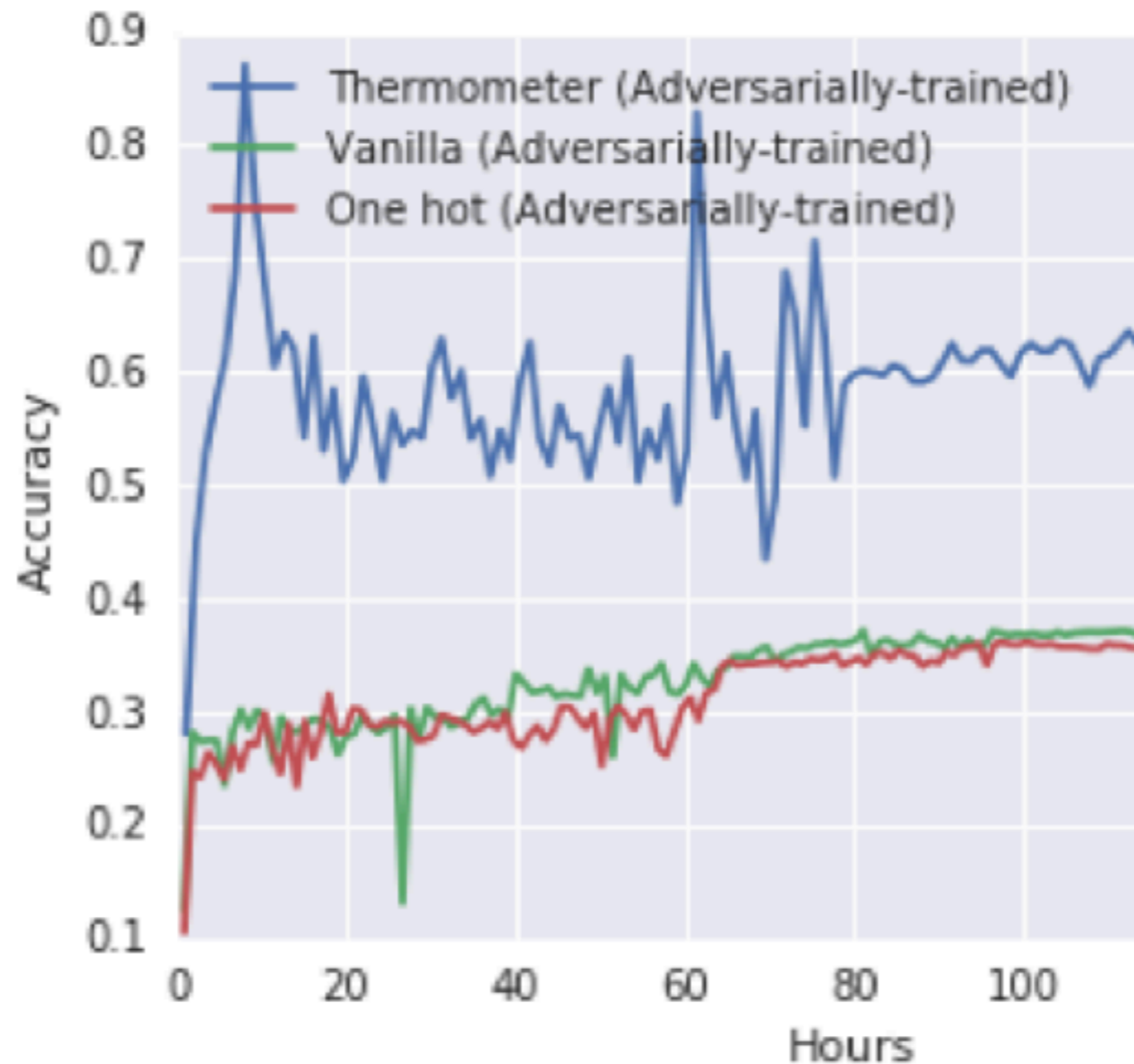
[0000000001]

Observation: PixelRNN shows one-hot codes work

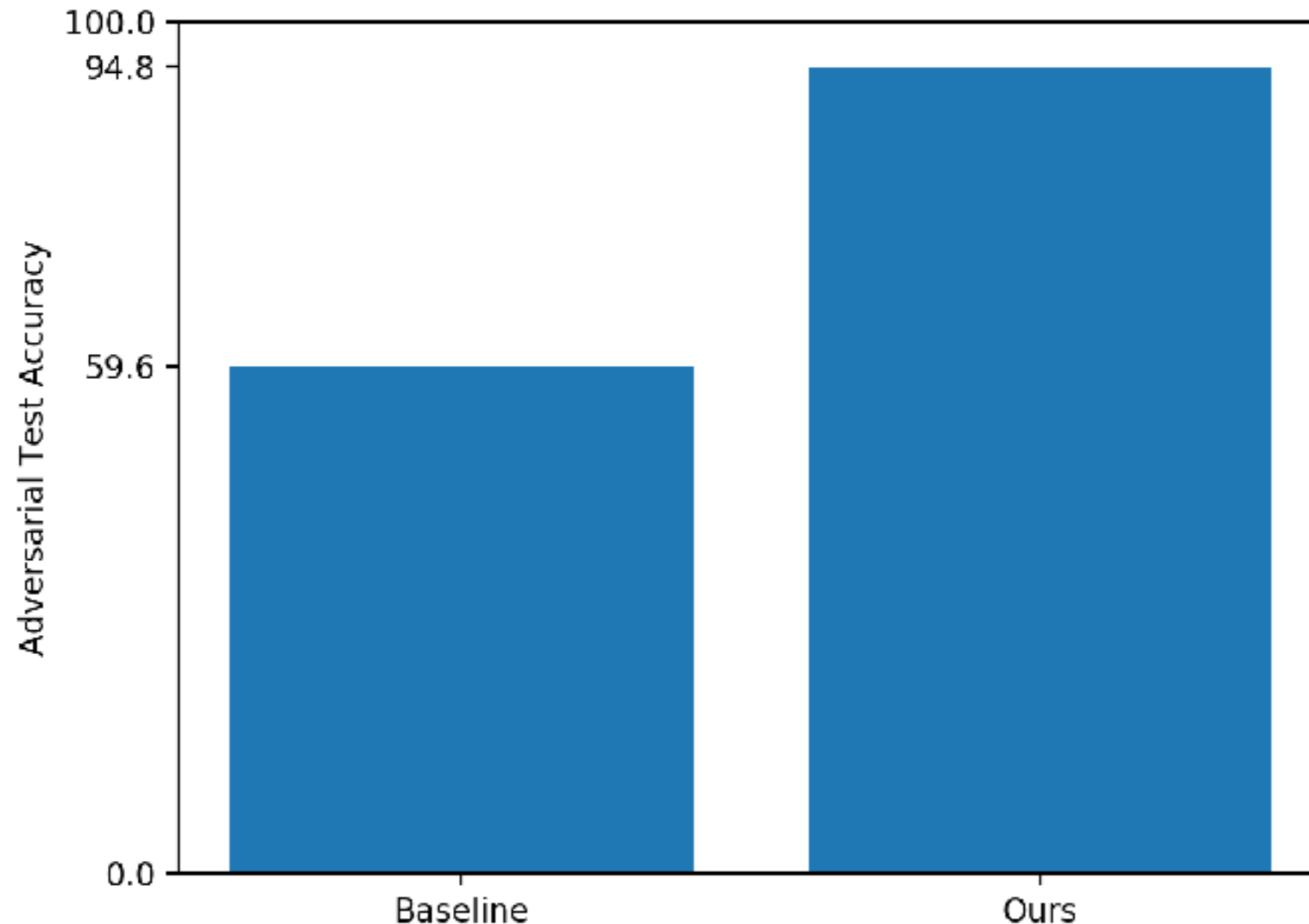




Fast Improvement Early in Learning

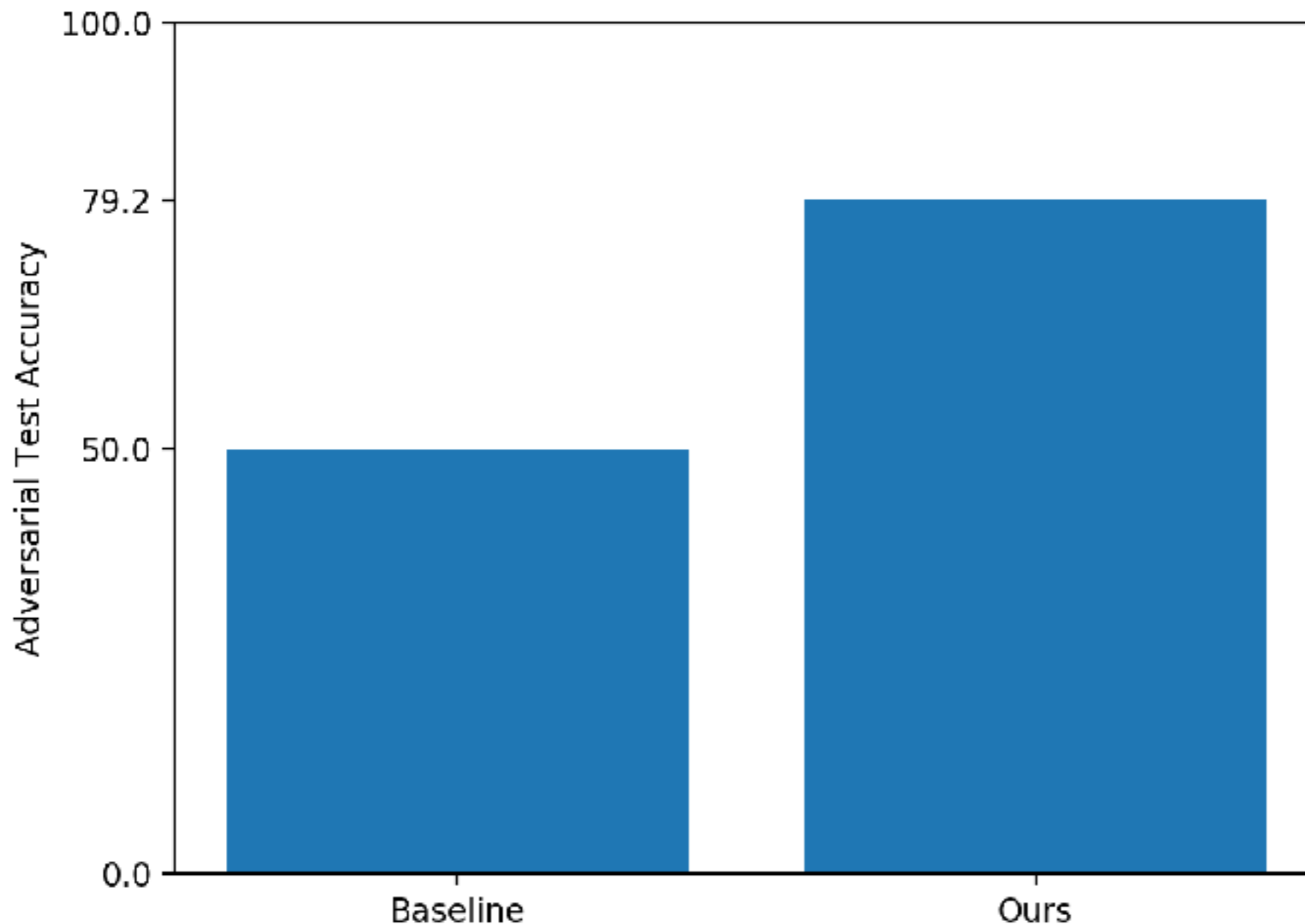


Large improvements on SVHN white box attacks



5 years ago,
this would have
been SOTA
on *clean* data

Large Improvements against CIFAR-10 white box attacks



6 years ago,
this would have
been SOTA
on *clean* data

Other results

- Improvement on CIFAR-100
 - (Still very broken)
- Improvement on MNIST
 - Please quit caring about MNIST

Caveats

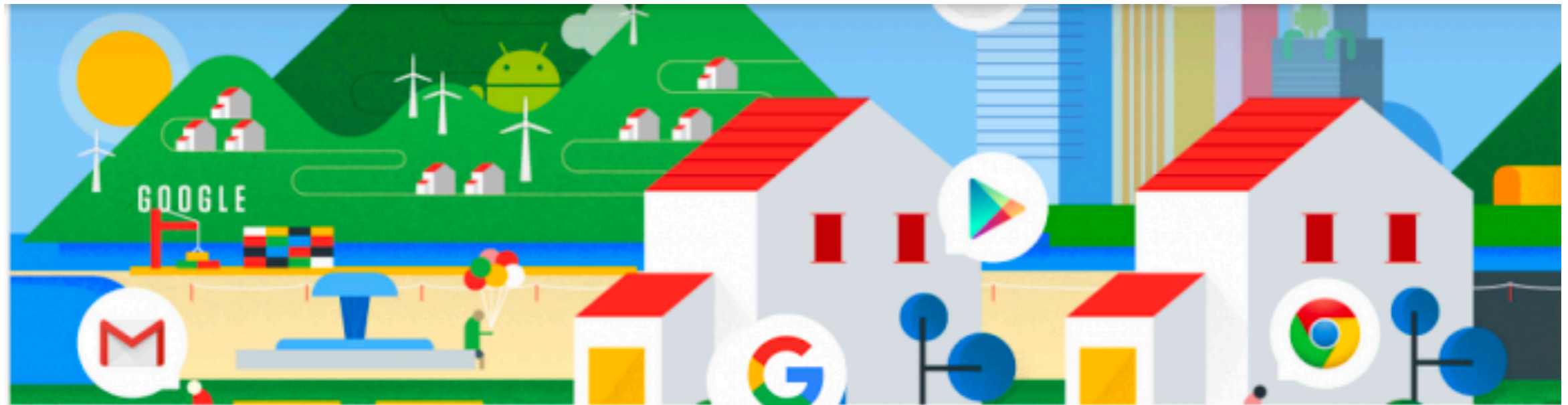
- Slight drop in accuracy on clean examples
- Only small improvement on black-box adversarial examples

Get involved!

<https://github.com/tensorflow/cleverhans>



g.co/airesidency



Google AI Resident, 2018 Start (Fixed-Term Employee)

Google

Software Engineering

Mountain View, CA, USA

[APPLY](#)

