

Understanding Deep Learning Requires Rethinking Generalization

Chiyuan Zhang
CSAIL, CBMM, MIT

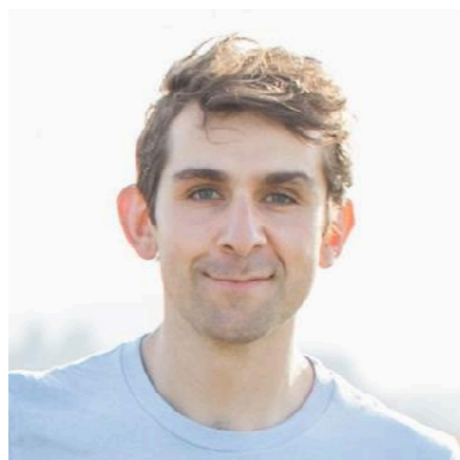
Poster: Wednesday Morning C23



Chiyuan Zhang



Samy Bengio



Moritz Hardt

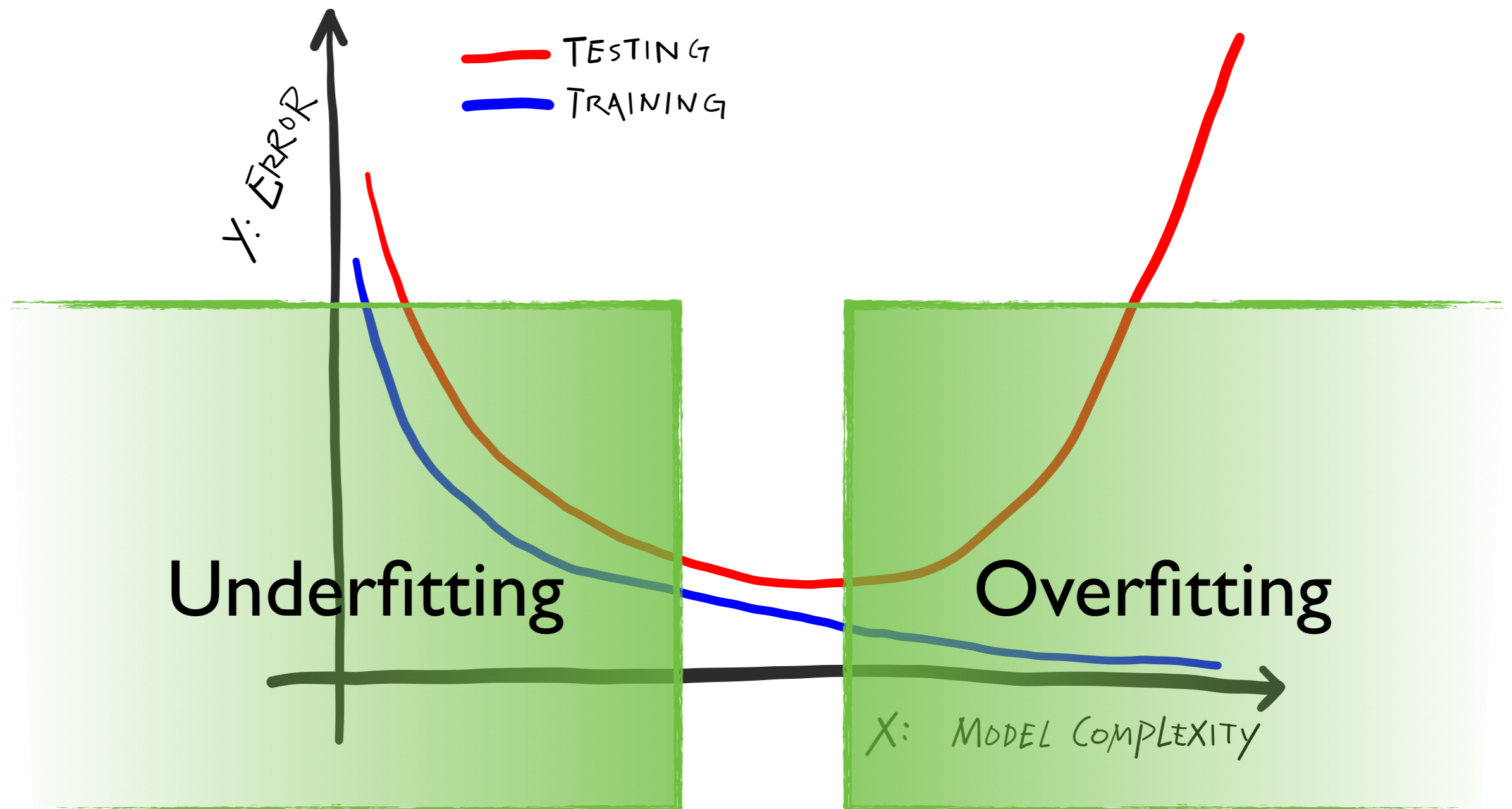


Benjamin Recht

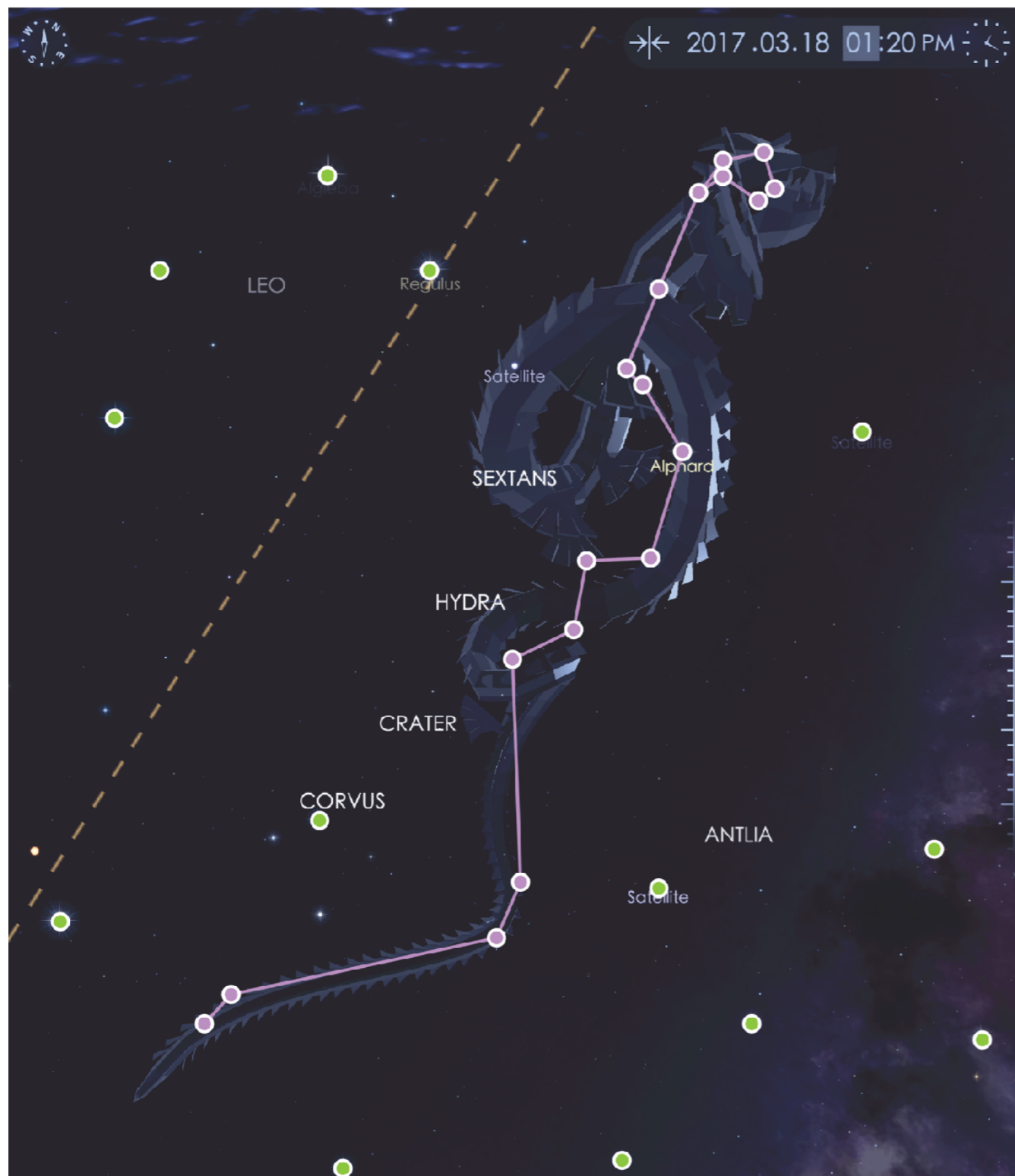


Oriol Vinyals

Model Selection



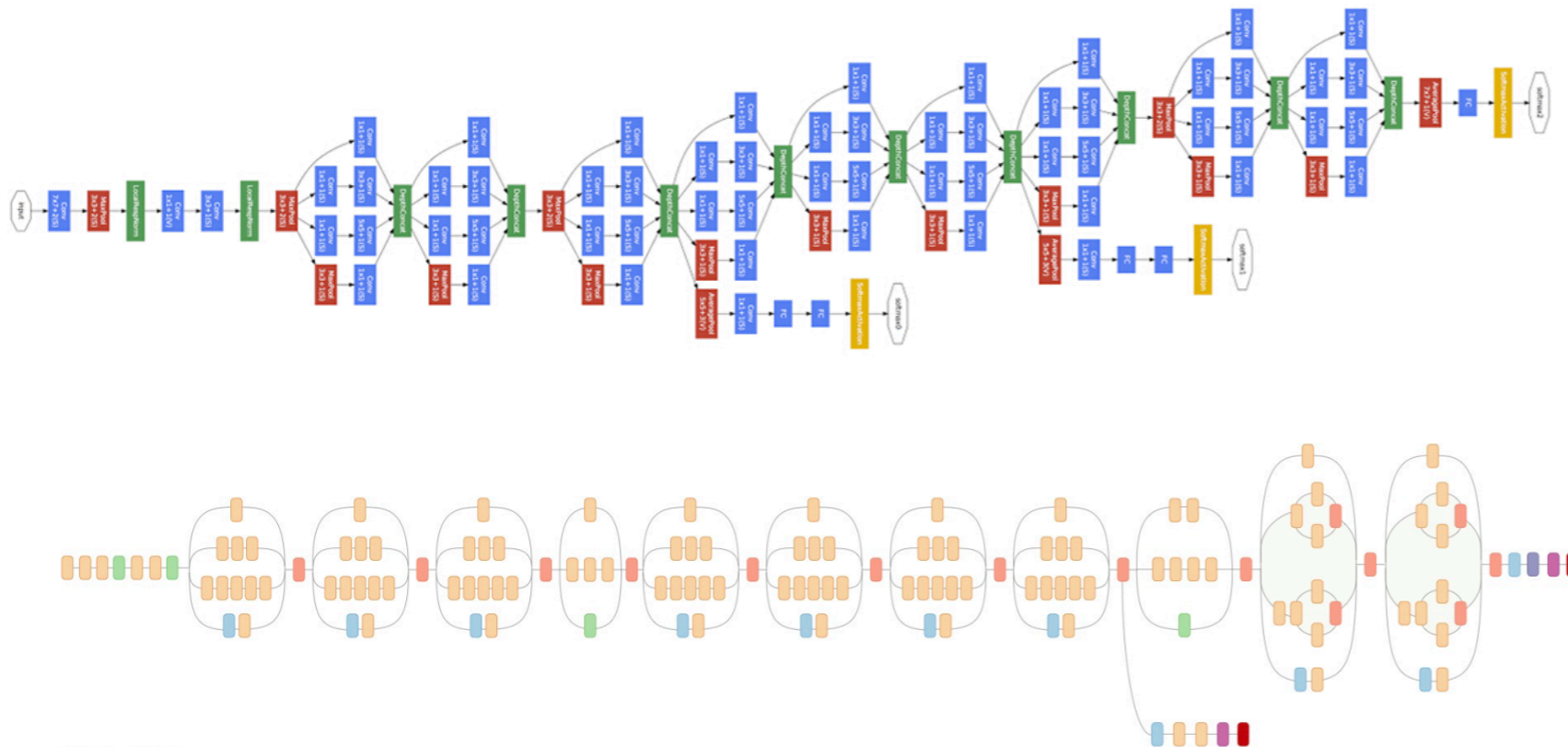
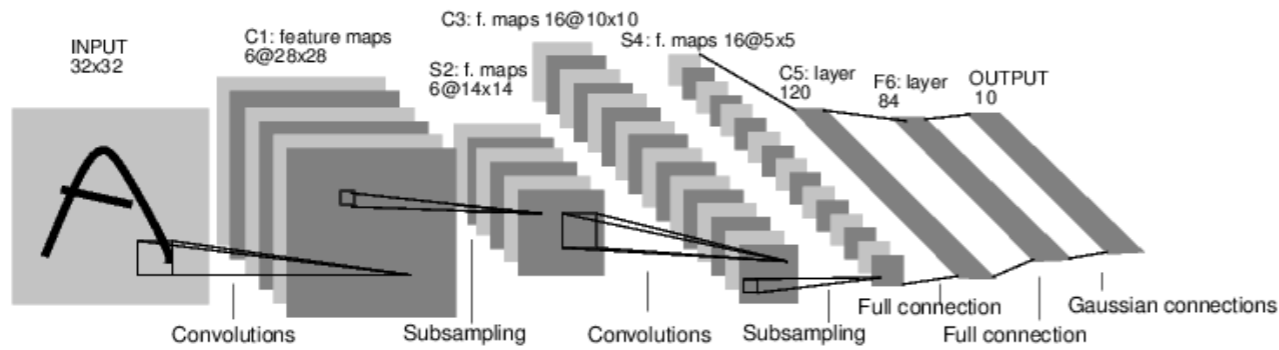
Over-parameterized Models



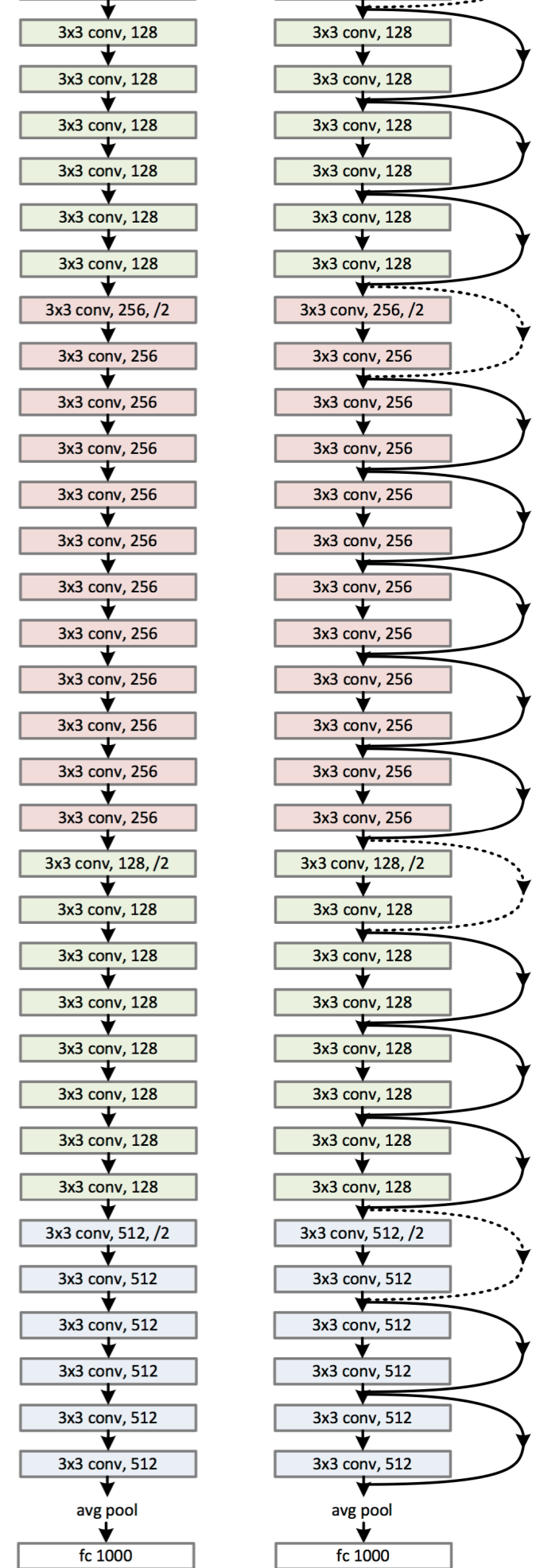
What are
the purple dots?

A Water Snake
The Constellation Hydra

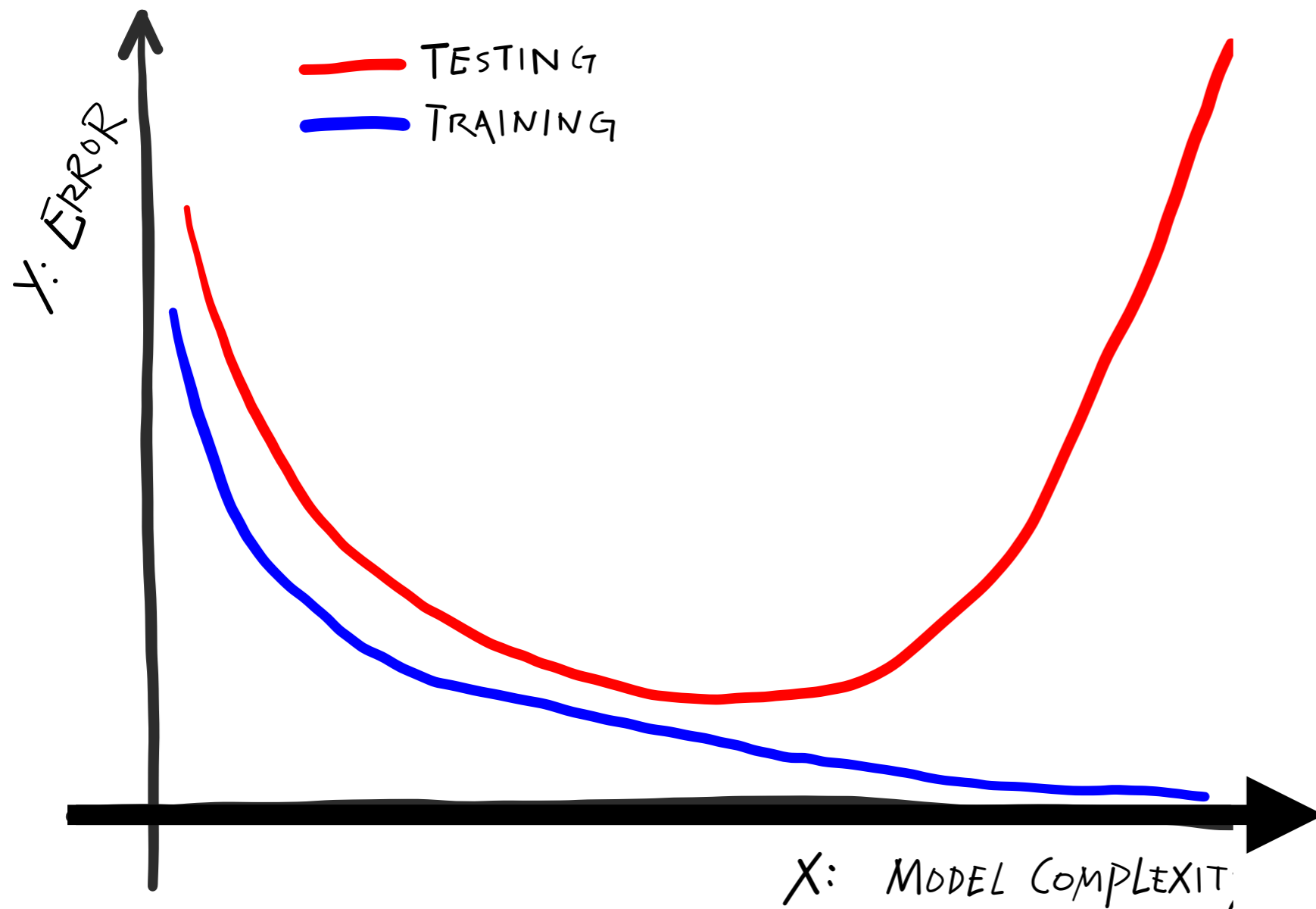
Deep Learning



- Convolution
- AvgPool
- MaxPool
- Concat
- Dropout
- Fully connected
- Softmax



Bias — Variance



Deep Learning

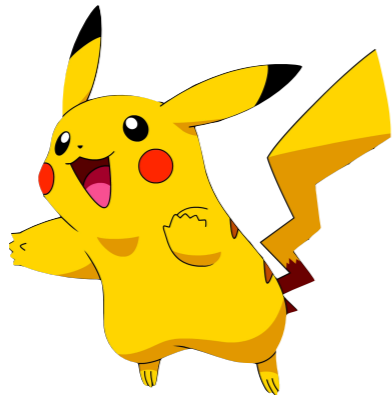


Parameter Count
Num Training Samples

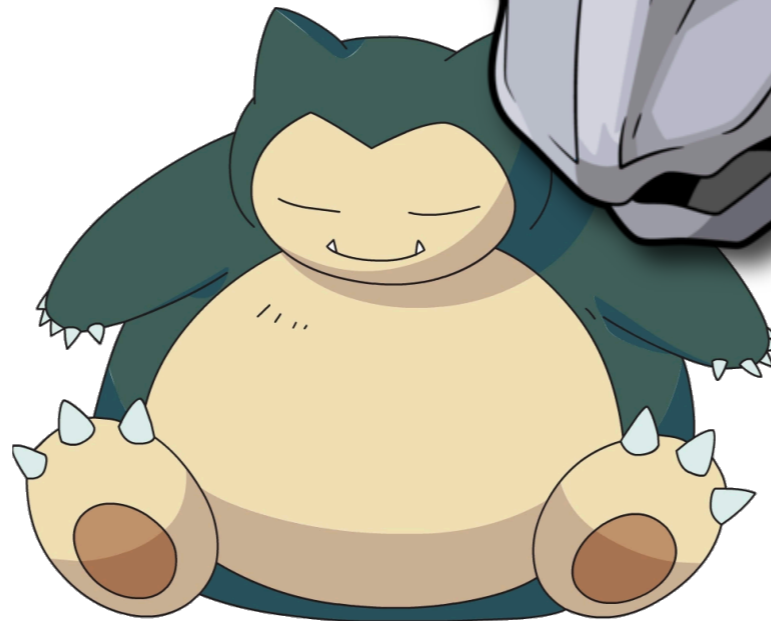
MLP 1x512
 $p/n: 24$



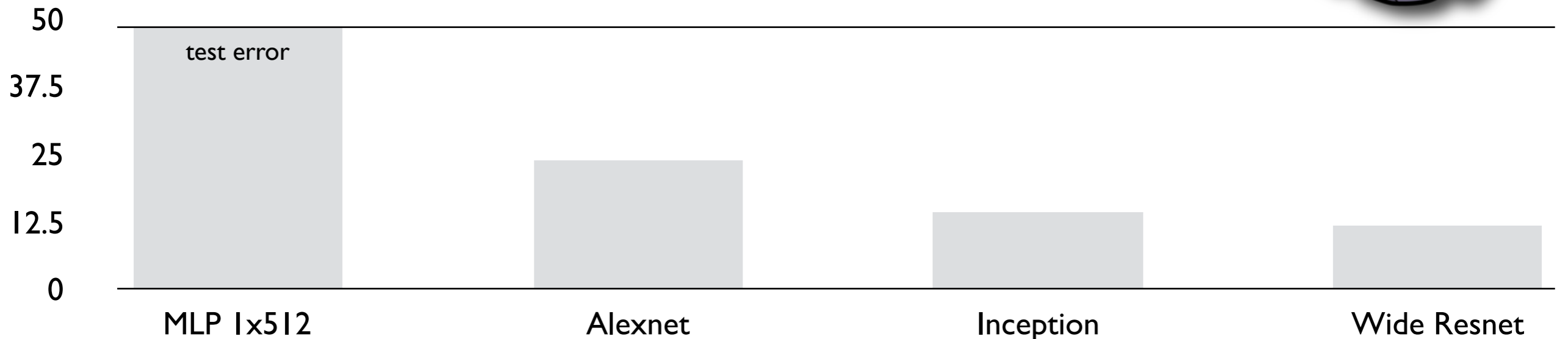
Alexnet
 $p/n: 28$



Inception
 $p/n: 33$



Wide Resnet
 $p/n: 179$



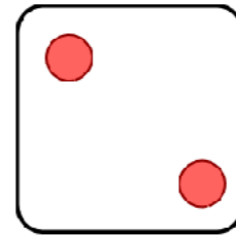
Randomization Test

Deep Neural Networks easily fit
random labels.

Random Label Dataset



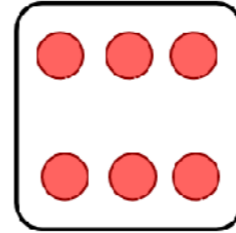
Dog



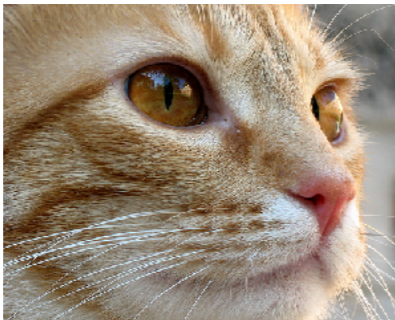
Cat



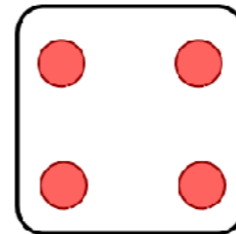
Flower



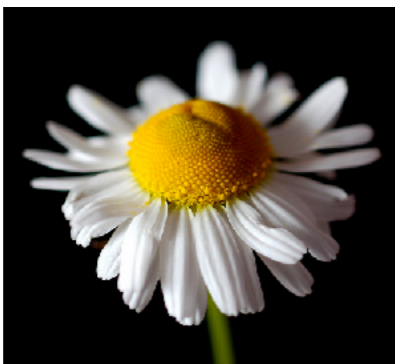
Dog



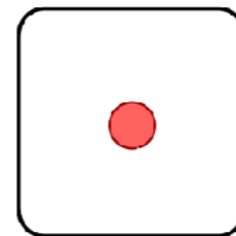
Cat



Bus



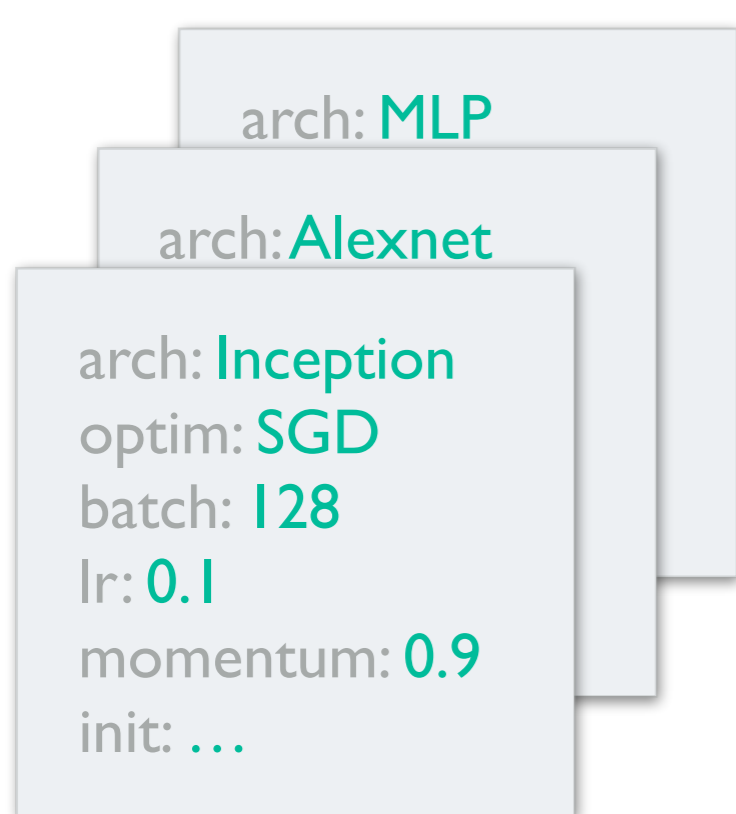
Flower



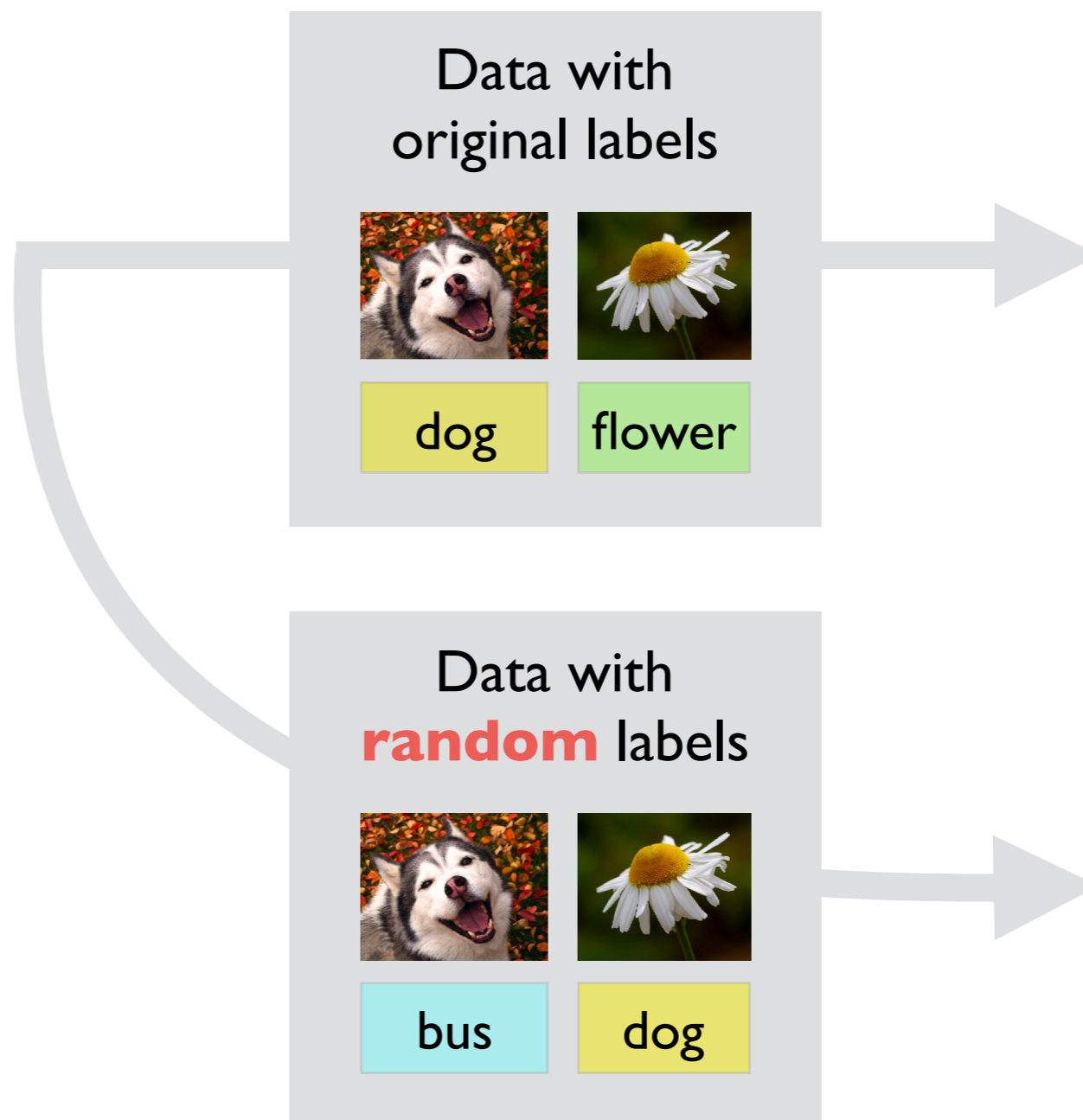
Bird

⋮

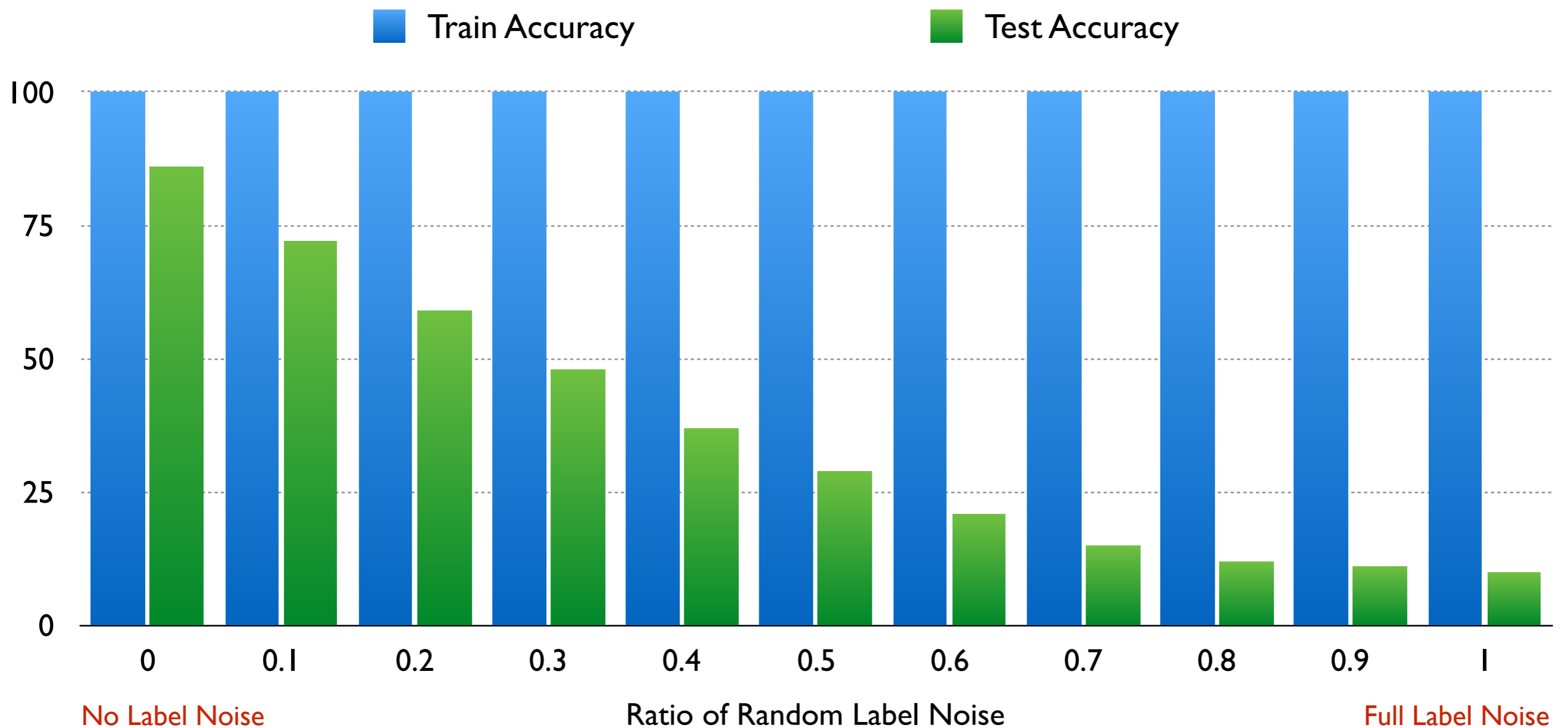
Randomization Test



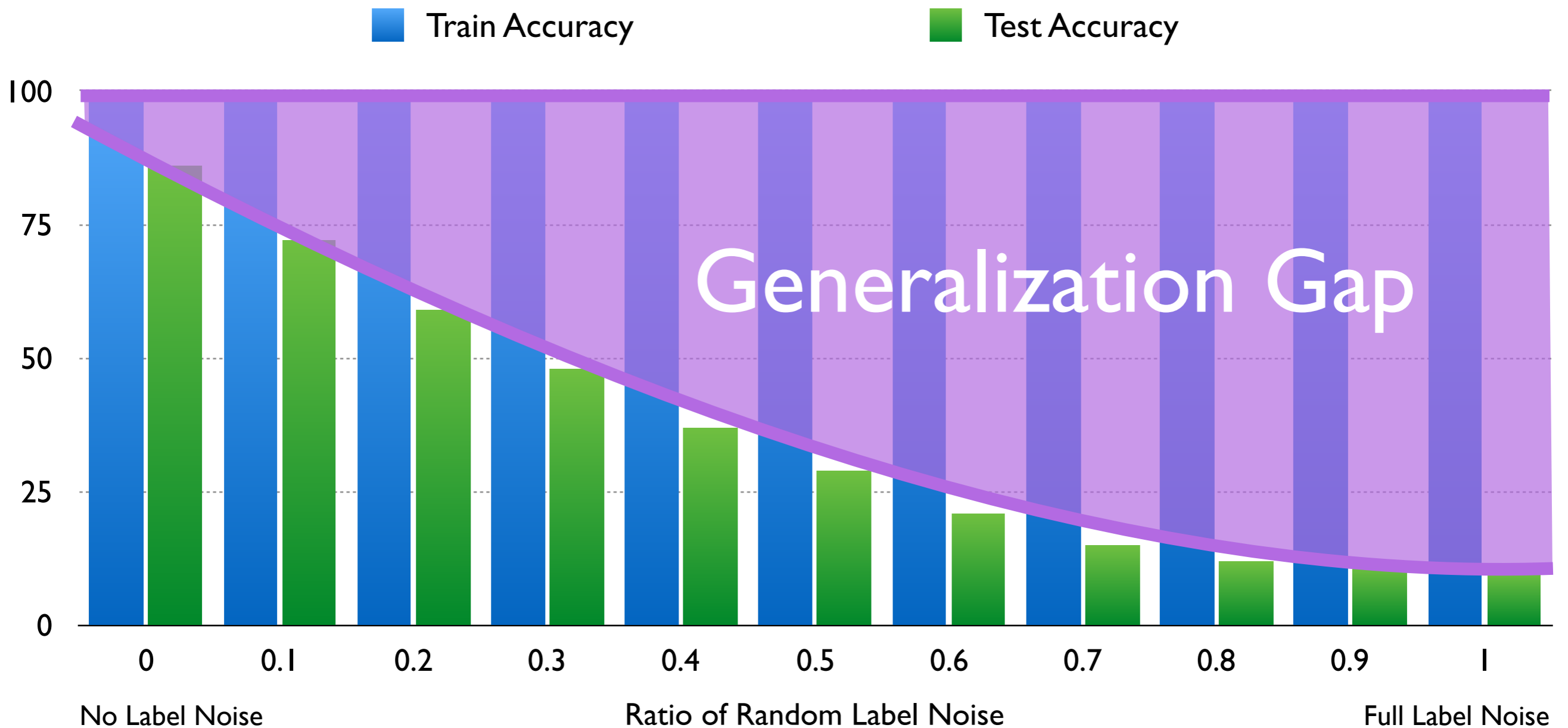
Recipes of Successful Models



Randomization Test



Randomization Test



Randomization Test

Deep Neural Networks easily fit
random labels.

Regularizers



⇐ Big Hypothesis Space

⇓ Regularized Models



Regularizers in Deep Learning

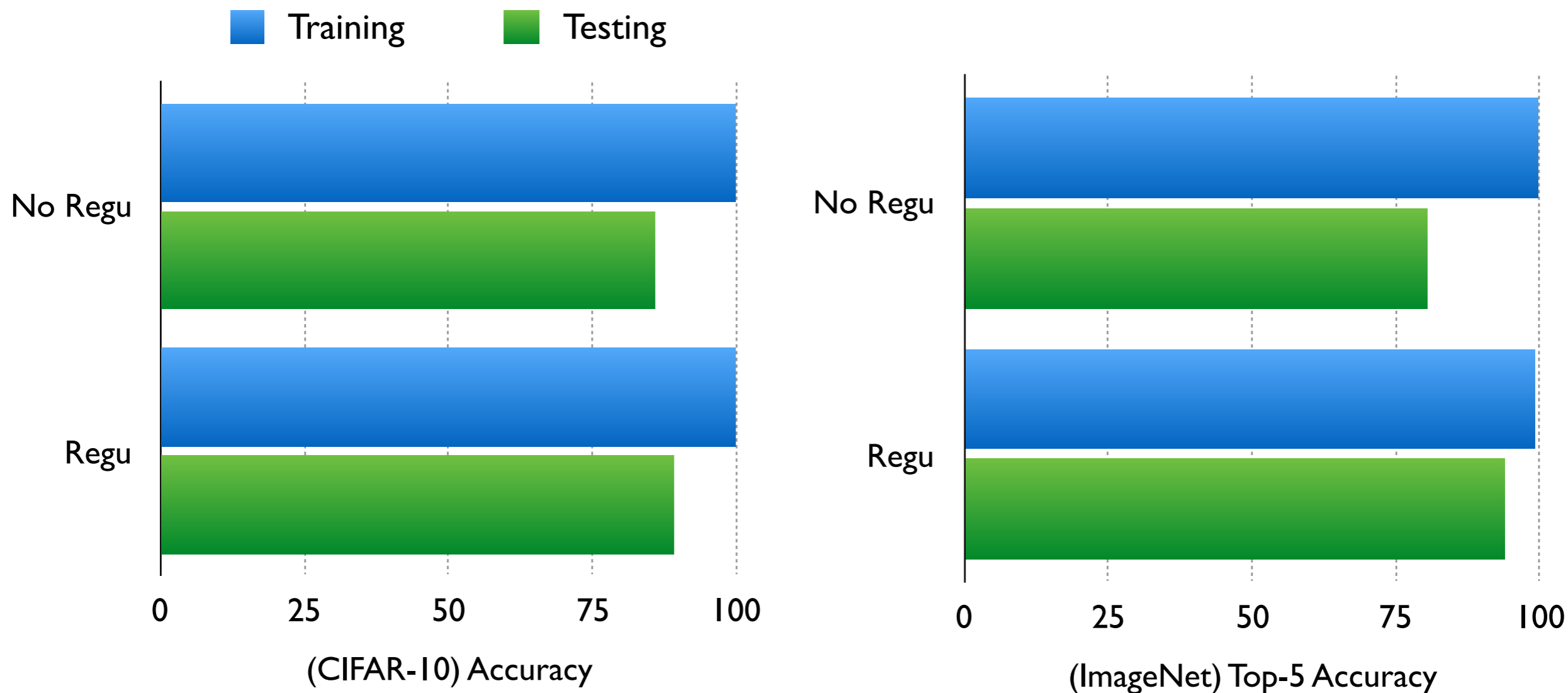
- Data augmentation: domain-specific transformations



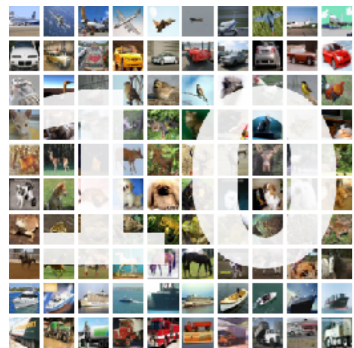
- Weight decay: l_2 -regularizer on weights
- Dropout*: randomly mask out responses



Fitting Natural Label with Regularizers



Fitting Random Label with Regularizers



Regularizer	Model	Training Accuracy
Weight decay	Inception	100%
	Alexnet	Failed to converge
	MLP 1x512	99.21%
Crop Augmentation*	Inception	99.93%



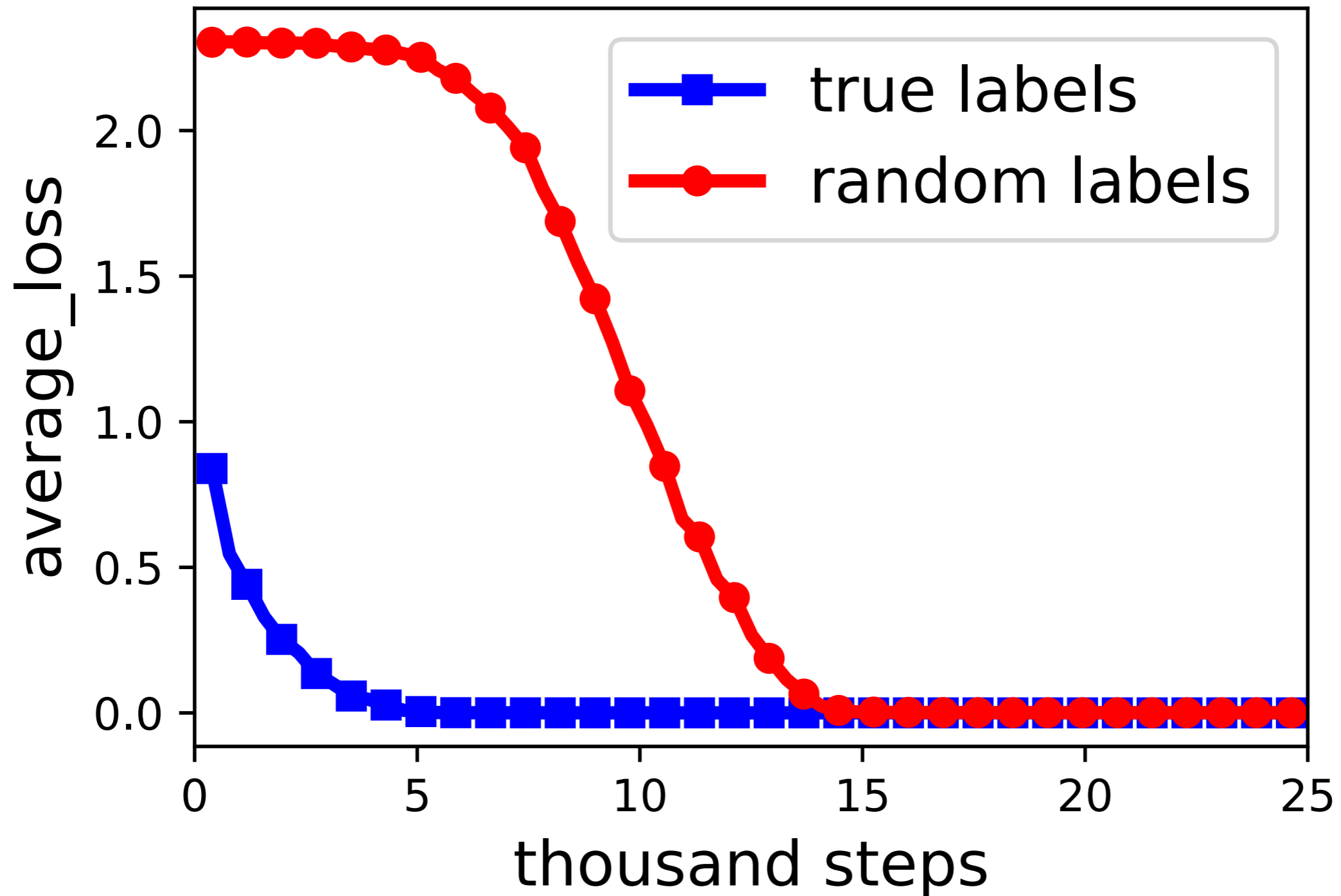
Regularizer	Model	Training top-5
Dropout	Inception V3	96.15%
Dropout + Weight decay		97.95%

*We need to tune the hyperparams a bit and run for more epochs for this to converge, see paper for details.

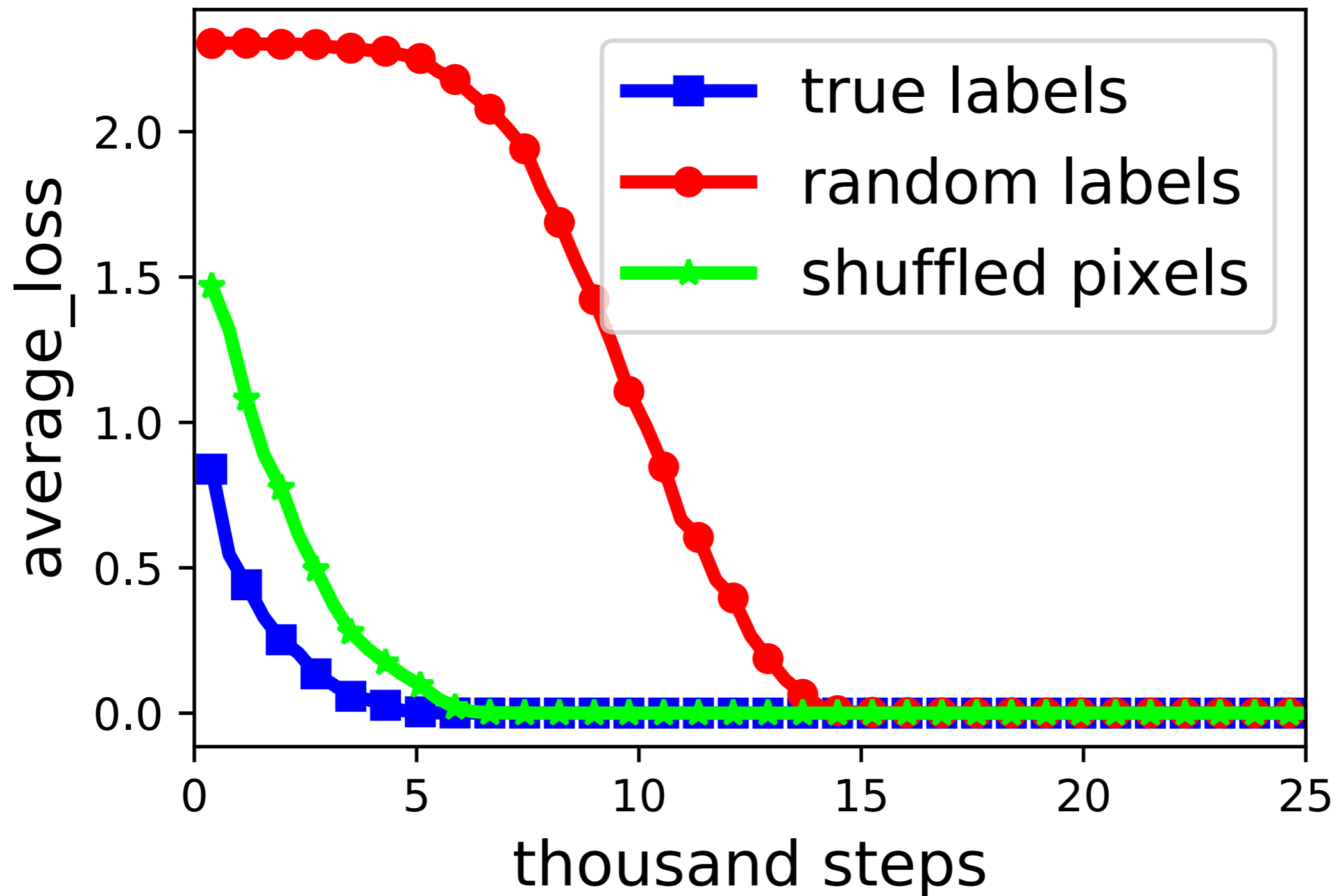
Implicit Regularization

A REGULARIZER is ^{anything} a ~~mechanism~~ that
~~constrain the model or empower the~~
~~data.~~ hurt the training Process.

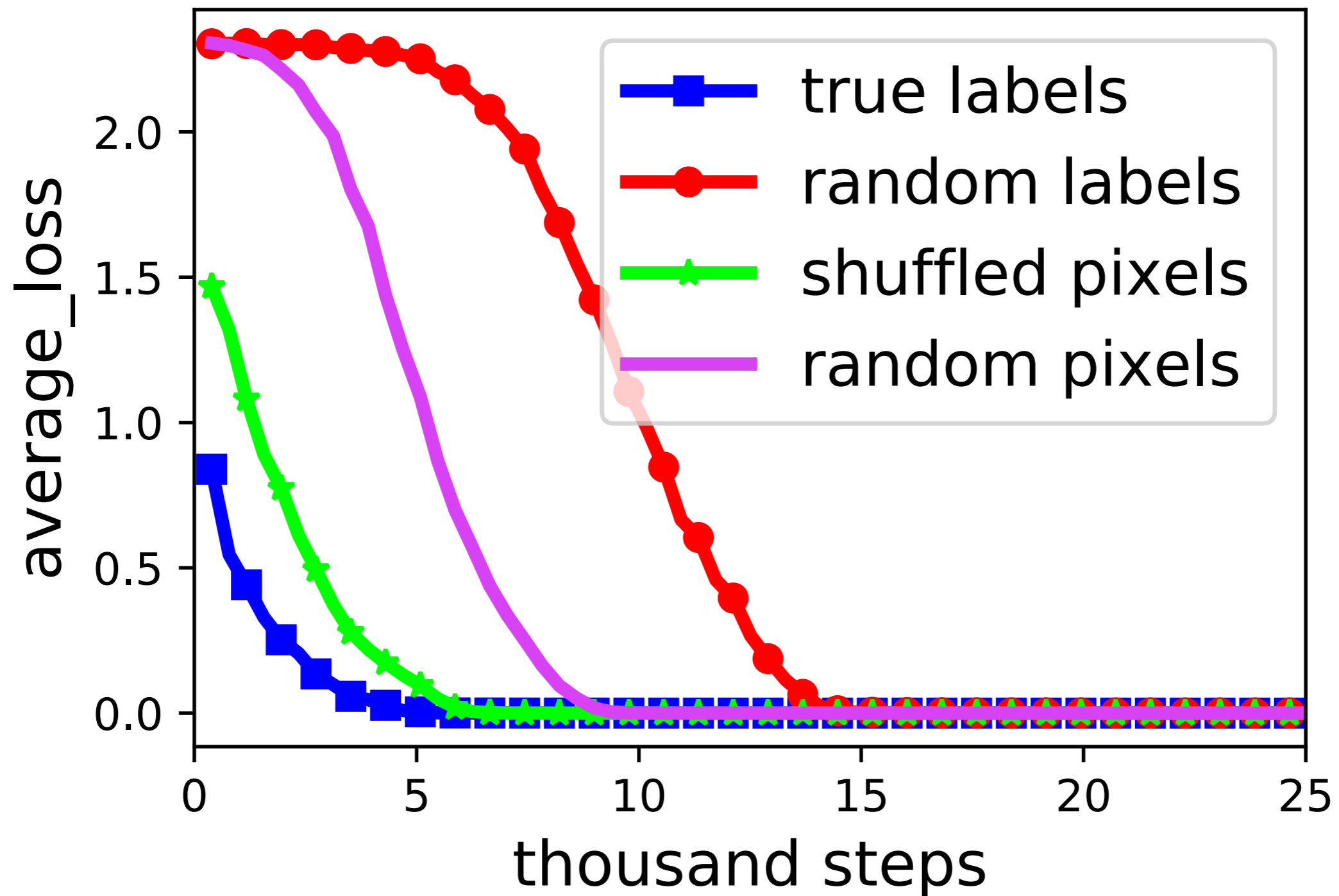
SGD fits Random Labels



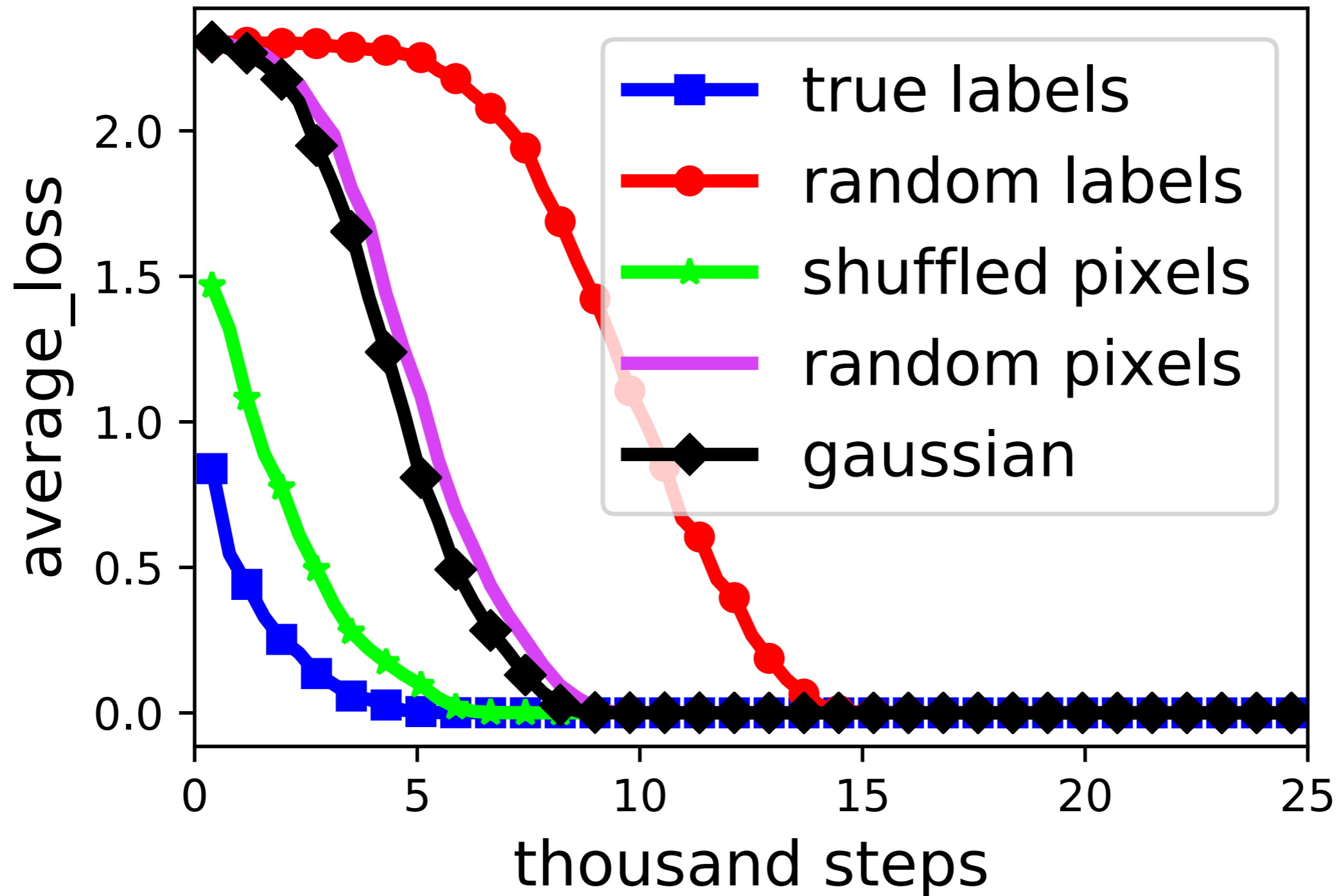
SGD fits Random Labels



SGD fits Random Labels



SGD fits Random Labels



SGD fits Random Labels

Optimization is easy for deep learning.

Conclusion

Simple experimental framework for understanding the **effective capacity** of deep learning models

Successful DeepNets are able to **shatter** the training set

Other formal measures of complexity for the models / algorithms / data distributions are needed to precisely explain the **over-parameterized regime**

Understanding Deep Learning Requires Rethinking Generalization  
 Chiyuan Zhang¹, Samy Bengio², Moritz Hardt², Benjamin Recht³, Oriol Vinyals⁴ ¹MIT, ²Google Brain, ³UC Berkeley, ⁴Google DeepMind

Introduction

Deep neural networks easily fit random labels. Explicit regularization may improve generalization performance, but is neither necessary nor by itself sufficient for controlling generalization error.

Effective Capacity via Randomization Tests

Motivation: # params ≠ complexity

Implications: Rademacher Complexity & VC-dimension

Role of Regularization (cont)

Deep neural networks shatter the training set.

With or without noise, with or without structure / pattern

Batch Normalization

Early Stopping

Analysis & Outlook

Finite Sample Expressivity: Capability to overfit

Linear Models & Implicit Regularization: SGD ⇒ Min-Norm

Conclusions

We presented a simple experimental framework for deriving and understanding a notion of effective capacity of machine learning models. We found:

- The effective capacity of several successful neural network architectures is large enough to shatter the training set.
- Optimization continues to be empirically easy even if the resulting model does not generalize.

Poster: Wednesday
 Morning (April 26th,
 10:30am to 12:30pm)
C23