

# On the Expressive Power of Deep Neural Networks

Maithra Raghu<sup>1 2</sup> Ben Poole<sup>3</sup> Jon Kleinberg<sup>1</sup> Surya Ganguli<sup>3</sup> Jascha Sohl Dickstein<sup>2</sup>

*Jiaru Zhang*  
5.21.2019

# Contents



**Introduction**

---



**Activation Pattern**

---



**Trajectory Length**

---



**Conclusion**

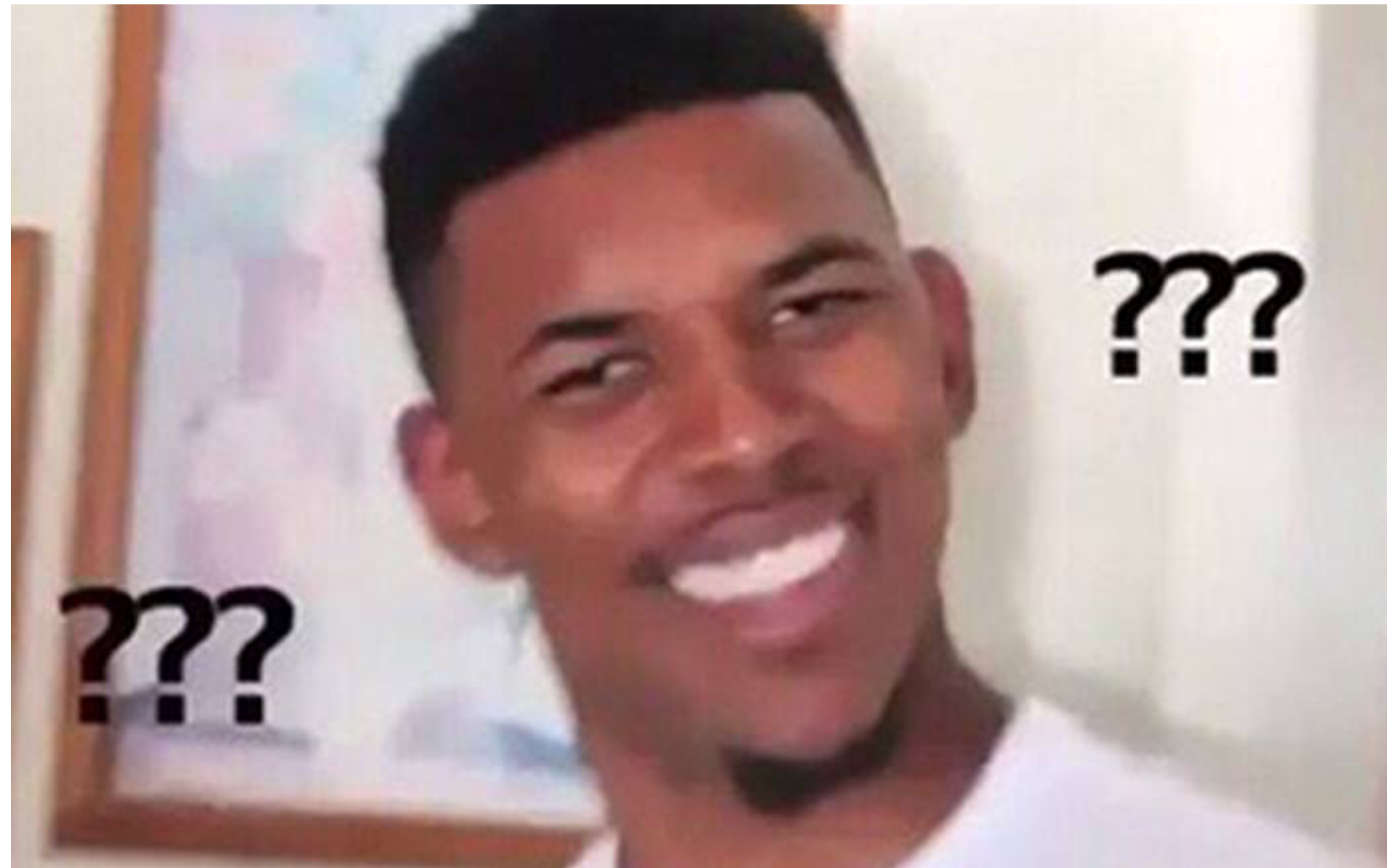
---

# Introduction

On the Expressive Power of Deep Neural Networks

# Introduction

## Expressive Power



1. What is it?
2. How to measure?
3. What determines it?
4. Usage?

# Contents



**Introduction**

---



**Activation Pattern**

---



**Trajectory Length**

---



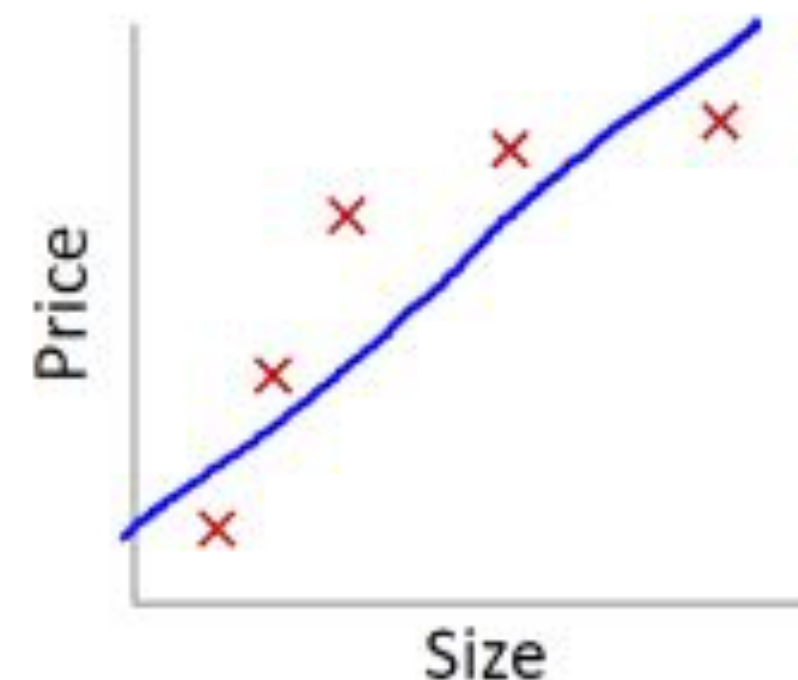
**Conclusion**

---

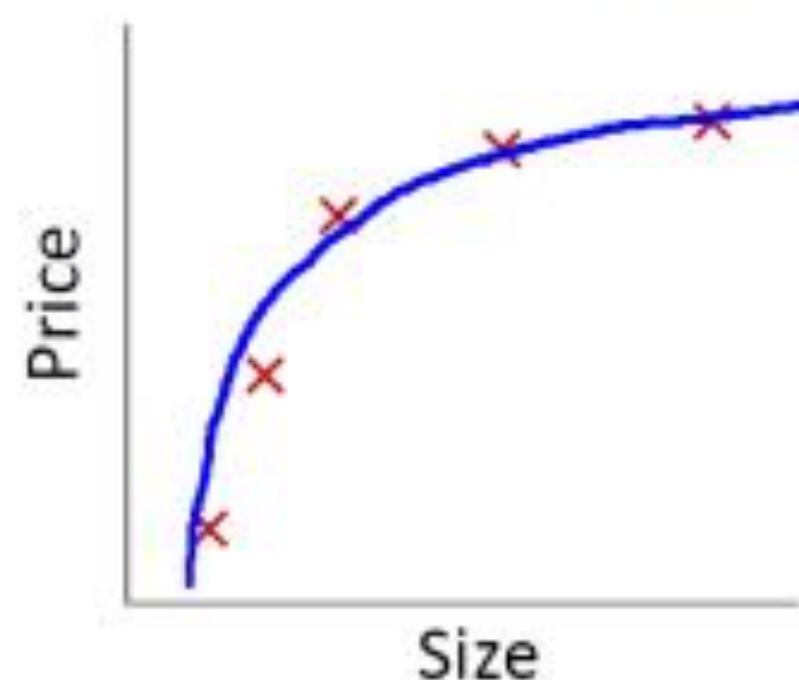
# Activation Pattern

What is Expressive Power (for a machine learning model)?

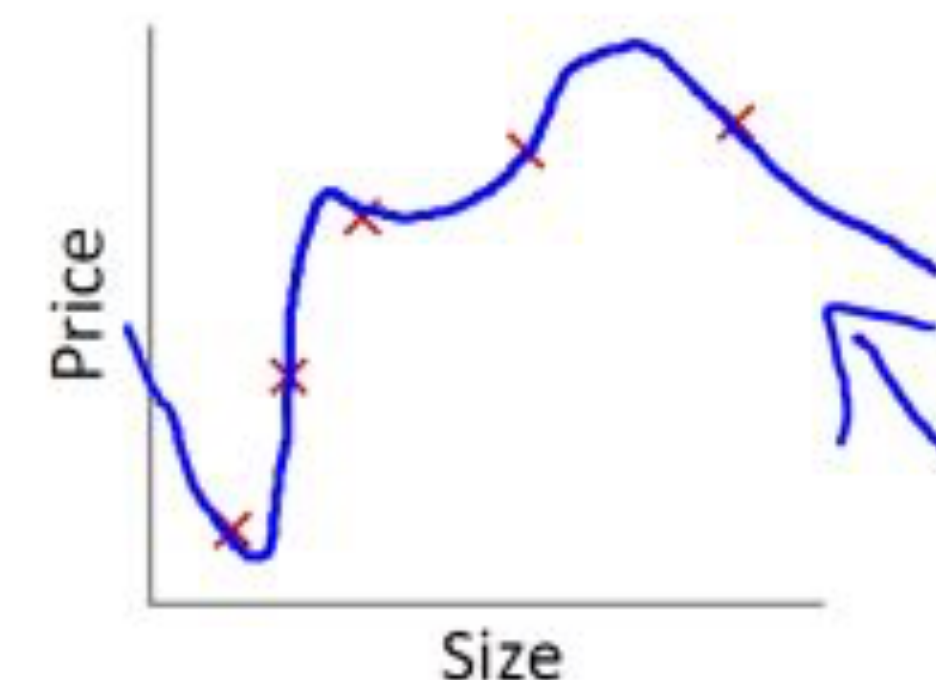
Example: Linear regression (housing prices)



$$\rightarrow \theta_0 + \theta_1 x$$



$$\rightarrow \theta_0 + \theta_1 x + \theta_2 x^2$$



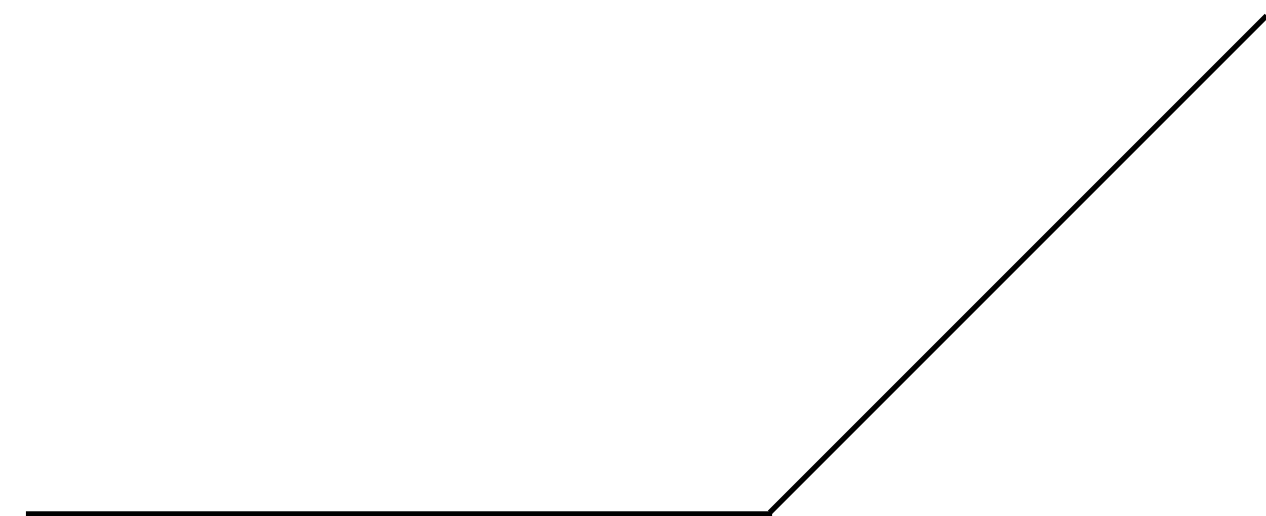
$$\rightarrow \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

Expressive power: Size of model space.

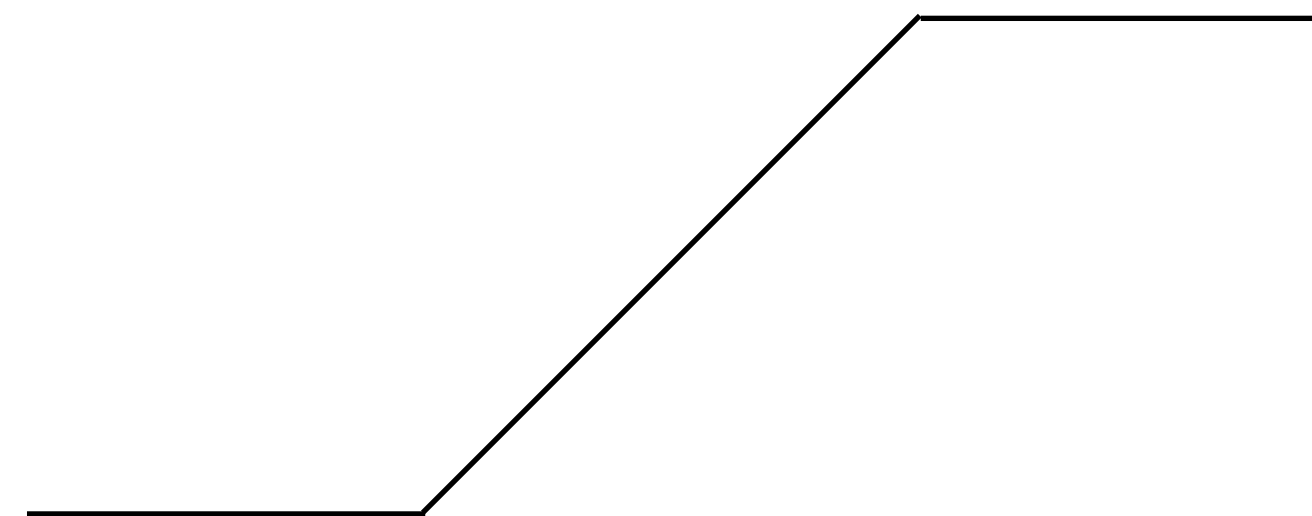
# Activation Pattern

## Settings for Neural Networks

- Only ReLU and hard tanh



ReLU

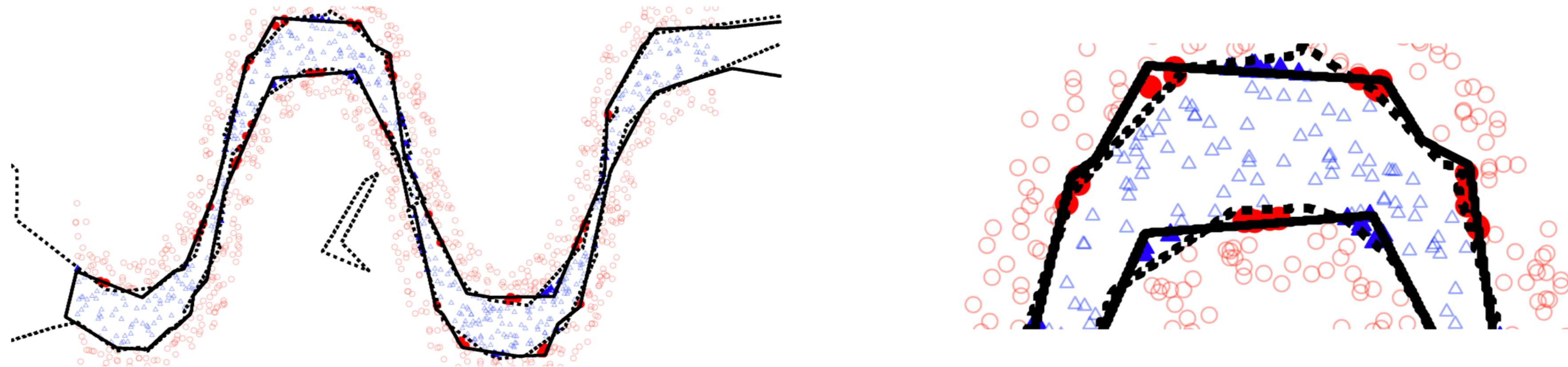


hard tanh

# Activation Pattern

## Settings for Neural Networks

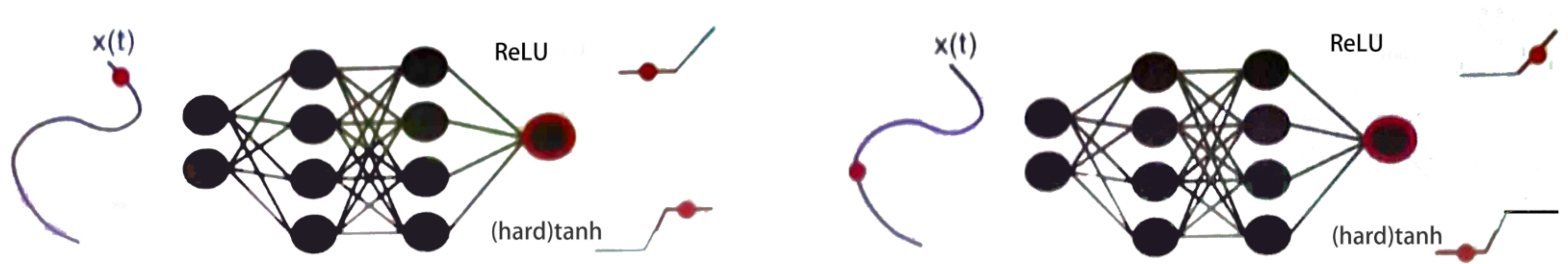
- Only ReLU and hard tanh





# Activation Pattern

How to measure?



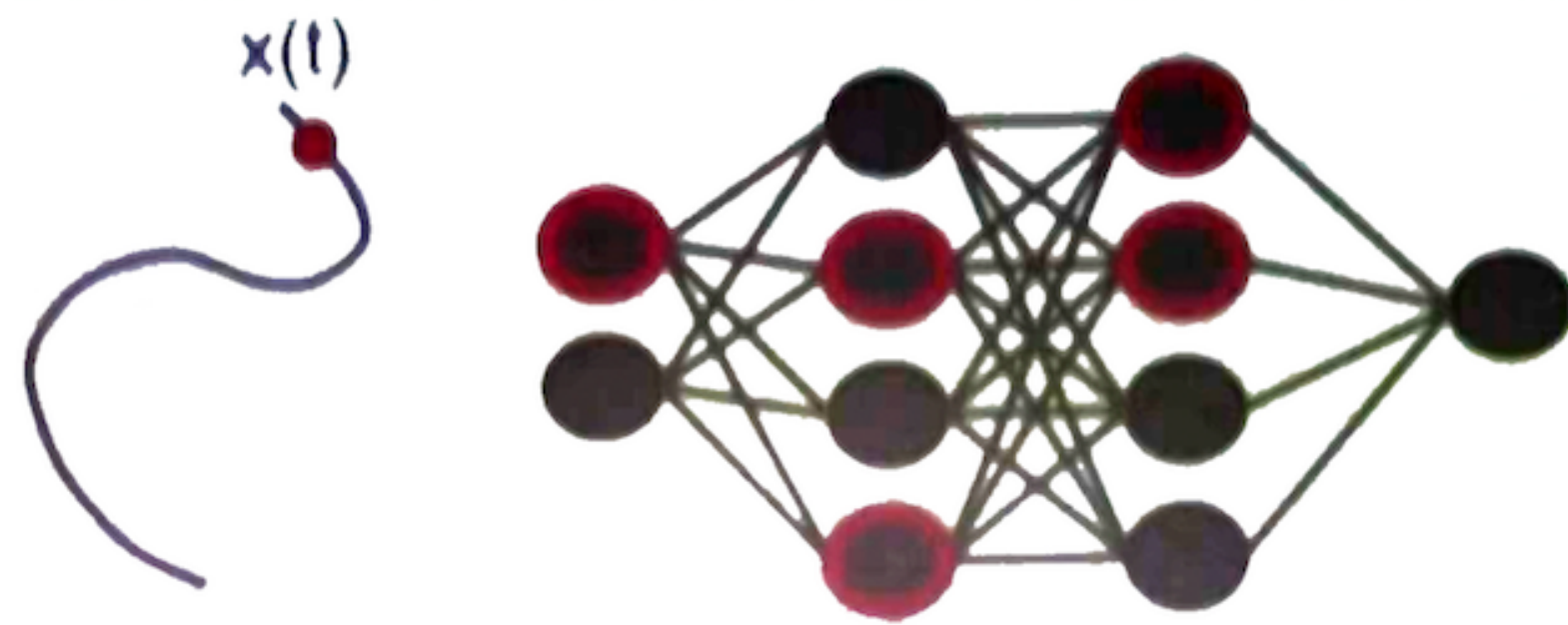
The number of linear regions

*Zaslavsky's theorem: a shallow network (i.e. one hidden layer), with the same number of parameters as a deep network, has a much smaller number of linear regions than the number achieved by their choice of weights  $W_0$  for the deep network.*

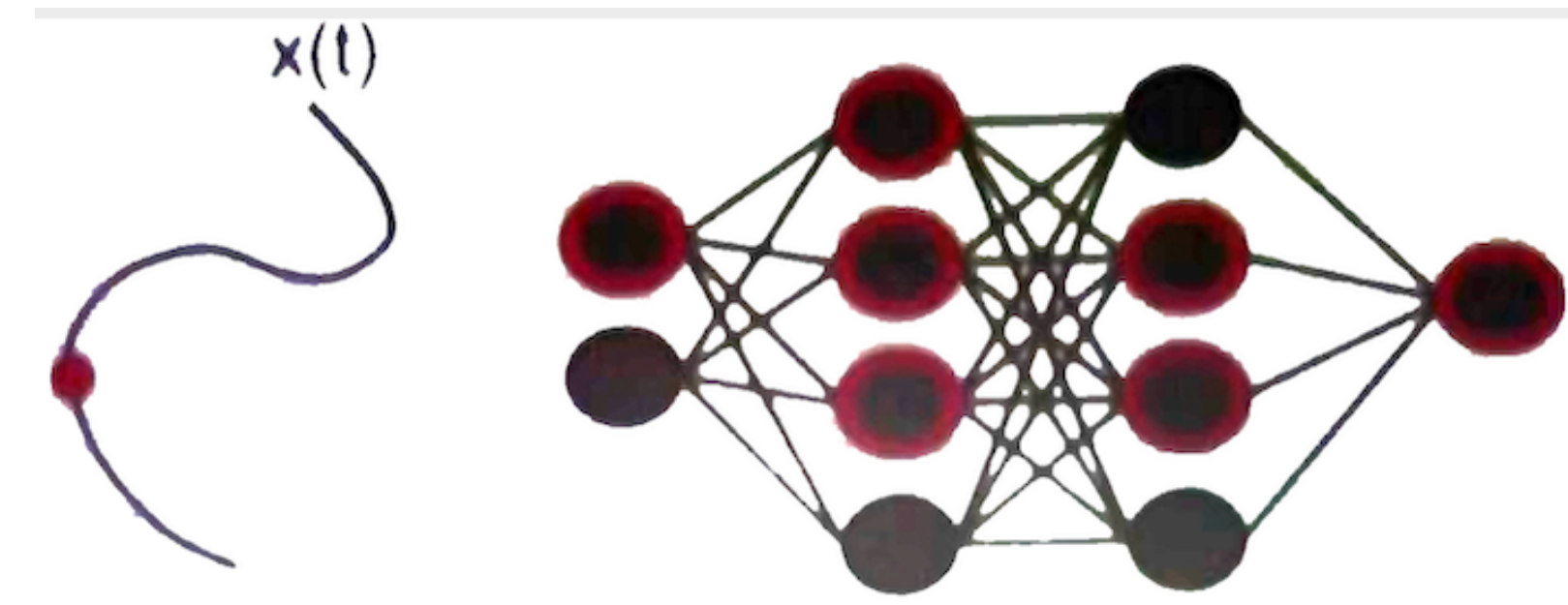
$$\forall W \quad \mathcal{T} \left( F_{A_1}([0,1]; W) \right) < \mathcal{T} \left( F_{A_1}([0,1]; W_0) \right)$$

# Activation Pattern

How to measure?



10010111000



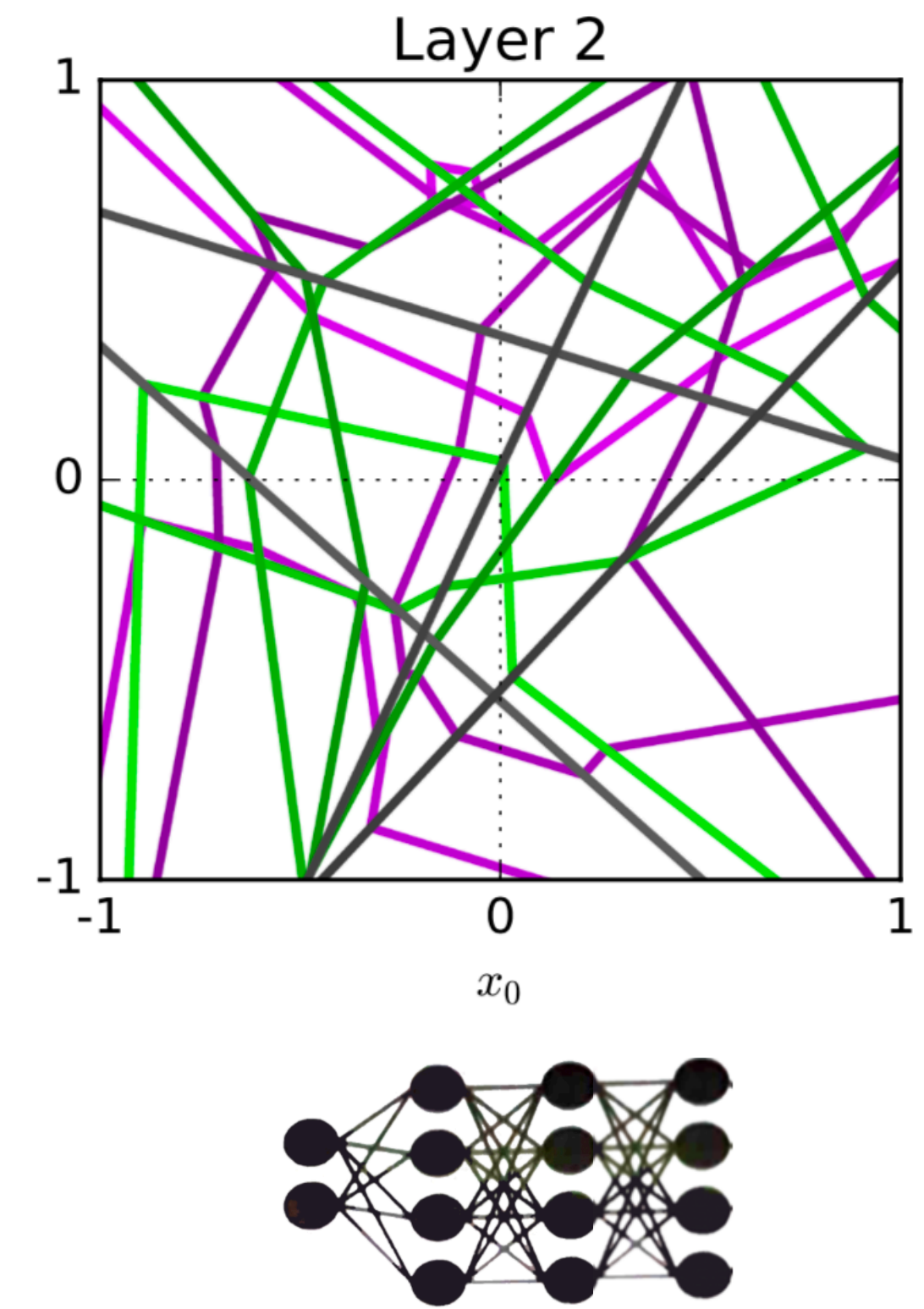
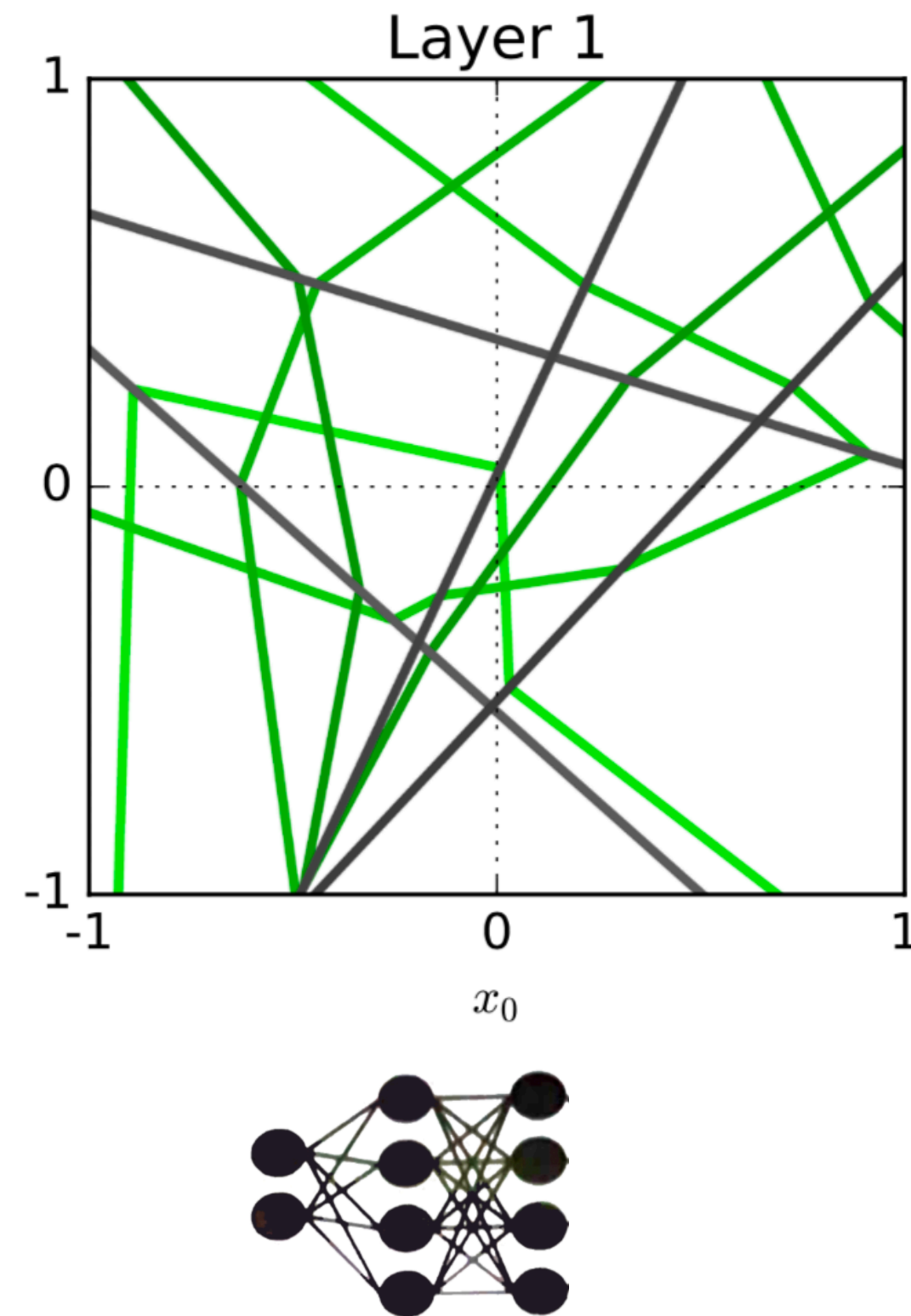
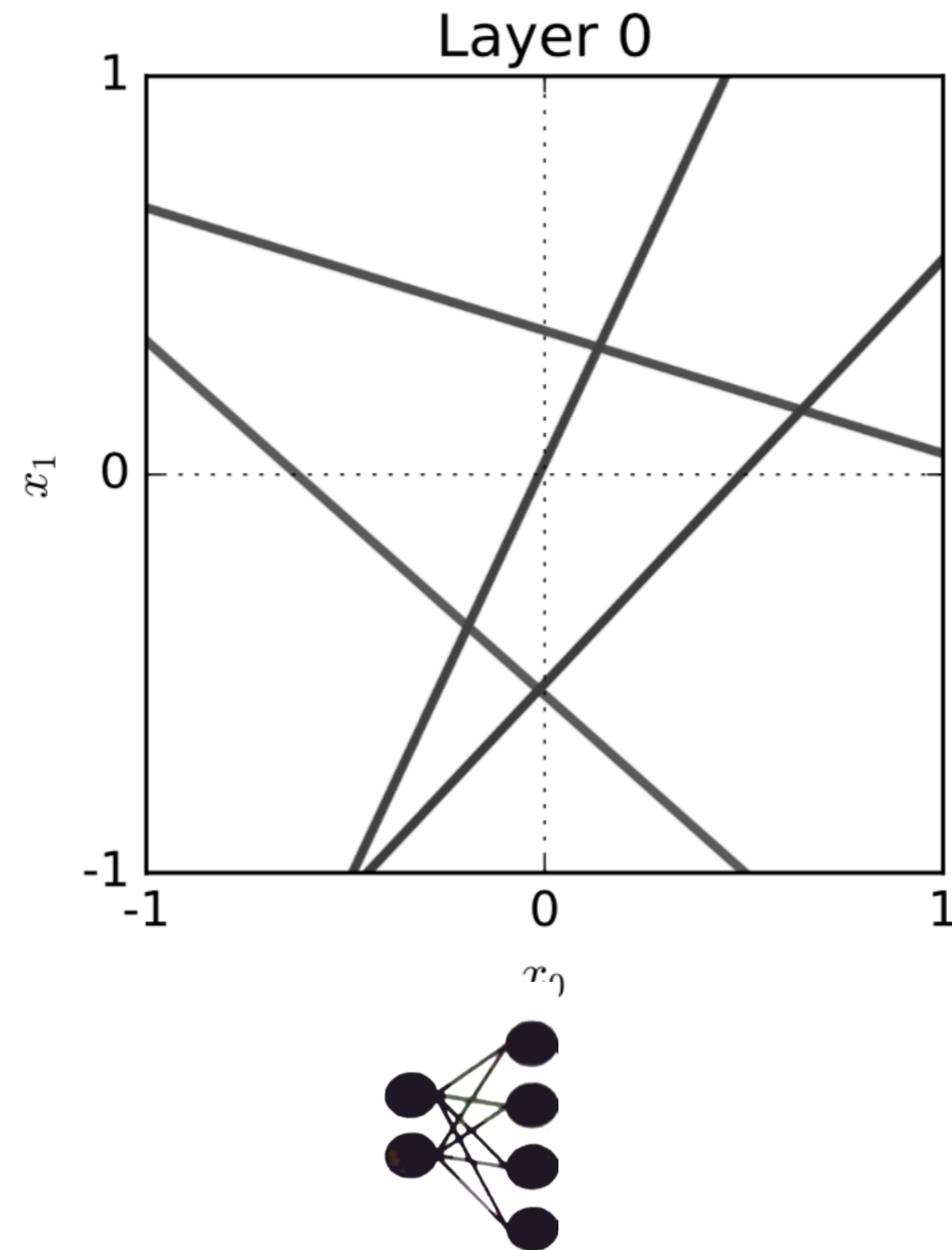
10111001101

number of Activation Patterns

# Activation Pattern

**Metric: number of Activation Patterns**

Activation patterns are in one-one correspondence with linear regions in input space.



# Activation Pattern

**What determines number of activation patterns?**

Upper bound grows linearly with depth and input dimension:

Given a network with:

- Depth  $n$
- Width  $k$
- Input dimension  $m$

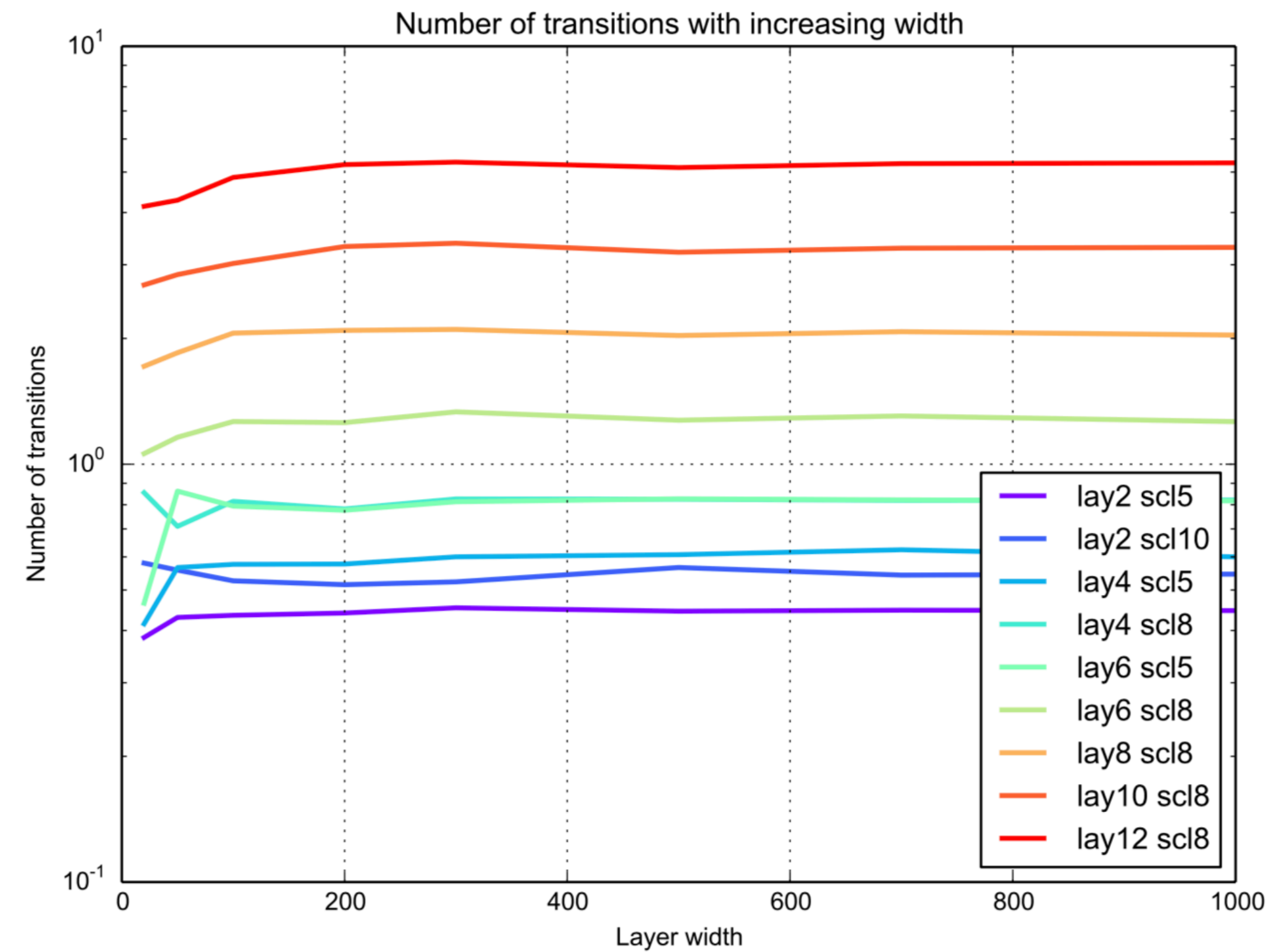
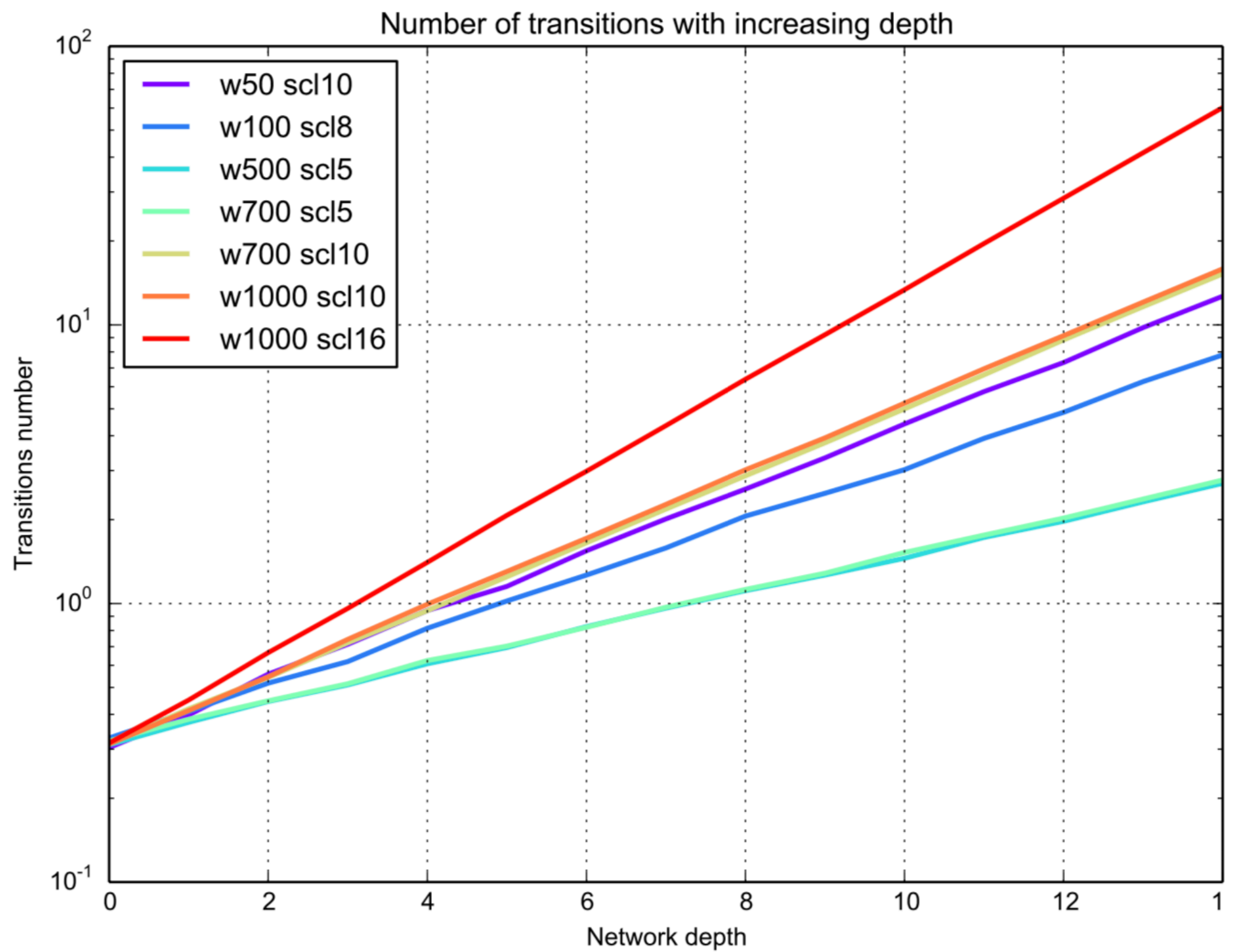
$$O(k^{mn}) \text{ (for ReLU)}$$

$$O((2k)^{mn}) \text{ (for hard tanh)}$$



# Activation Pattern

## Experiment Verification



# Activation Pattern

## Insight: A Derivation

Question: When number of total neurons are fixed, how to arrange them to get best expressive power?

Answer: Total  $N$  neurons,  $N/k$  depth,  $k$  width. Input dimension  $m$ .

$$\text{Maximize } O(k^{m \frac{N}{k}}) \quad \implies \quad k = e$$

Conclusion: When  $k \geq 3$ , APs decrease when  $k$  increases.

# Activation Pattern

## Review

1. What is it?
  2. How to measure?
  3. What determines it?
  4. Usage?
1. Number of states
  2. Activation Patterns, 10111001101
  3.  $O(k^{mn})$  (ReLU),  $O((2k)^{mn})$  (hard tanh)
  4.  $k = 3$

## Question?

# Contents



**Introduction**

---



**Activation Pattern**

---



**Trajectory Length**

---



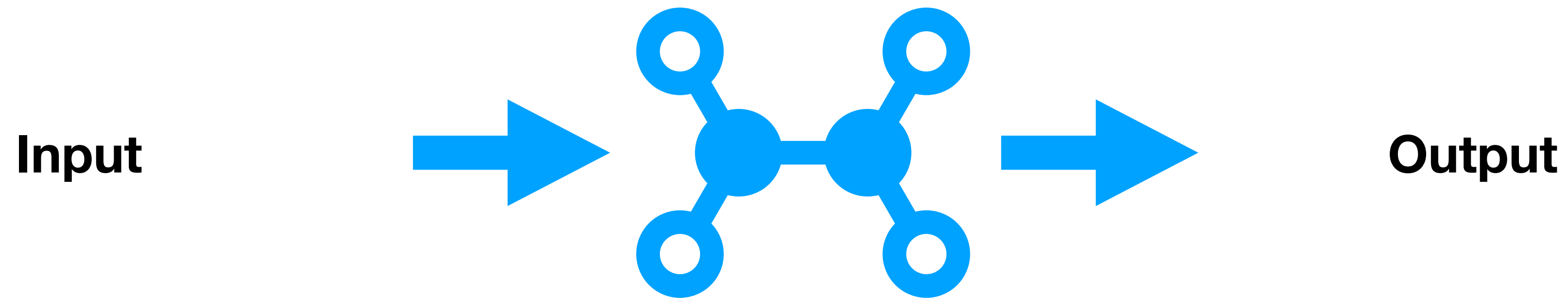
**Conclusion**

---



# Trajectory Length

**Rethink: What is Expressive Power?**



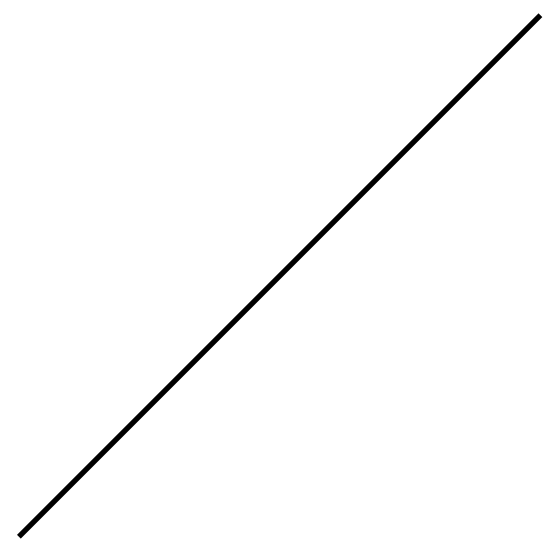
Neural network: A mapping from input to output.

Expressive Power: How complex the mapping is.

# Trajectory Length

Rethink: What is Expressive Power?

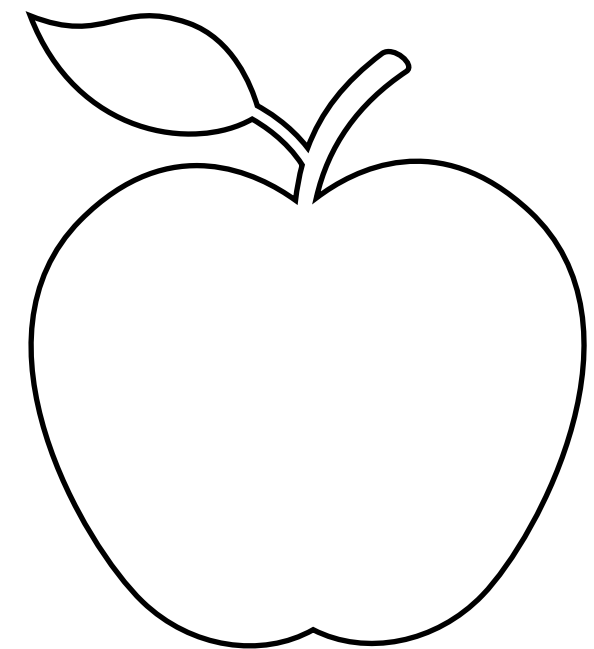
- Consider an (one-dimensional) input trajectory



$$x(t) = tx_1 + (1 - t)x_0$$



$$x(t) = \cos(\pi t/2)x_0 + \sin(\pi t/2)x_1$$

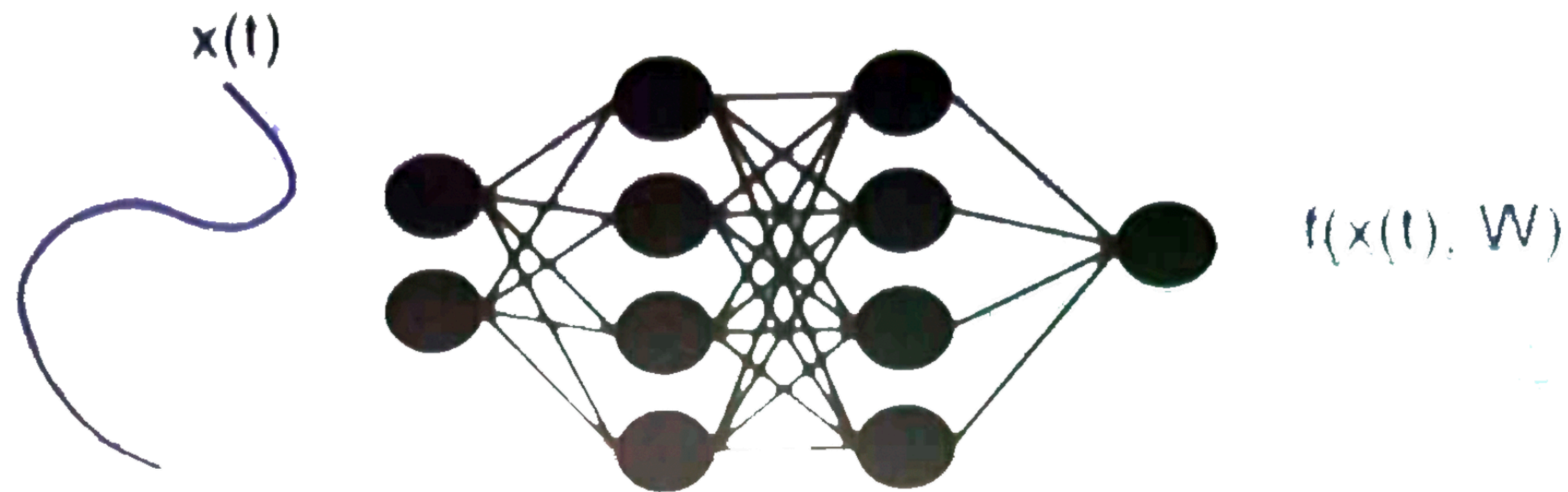


???

# Trajectory Length

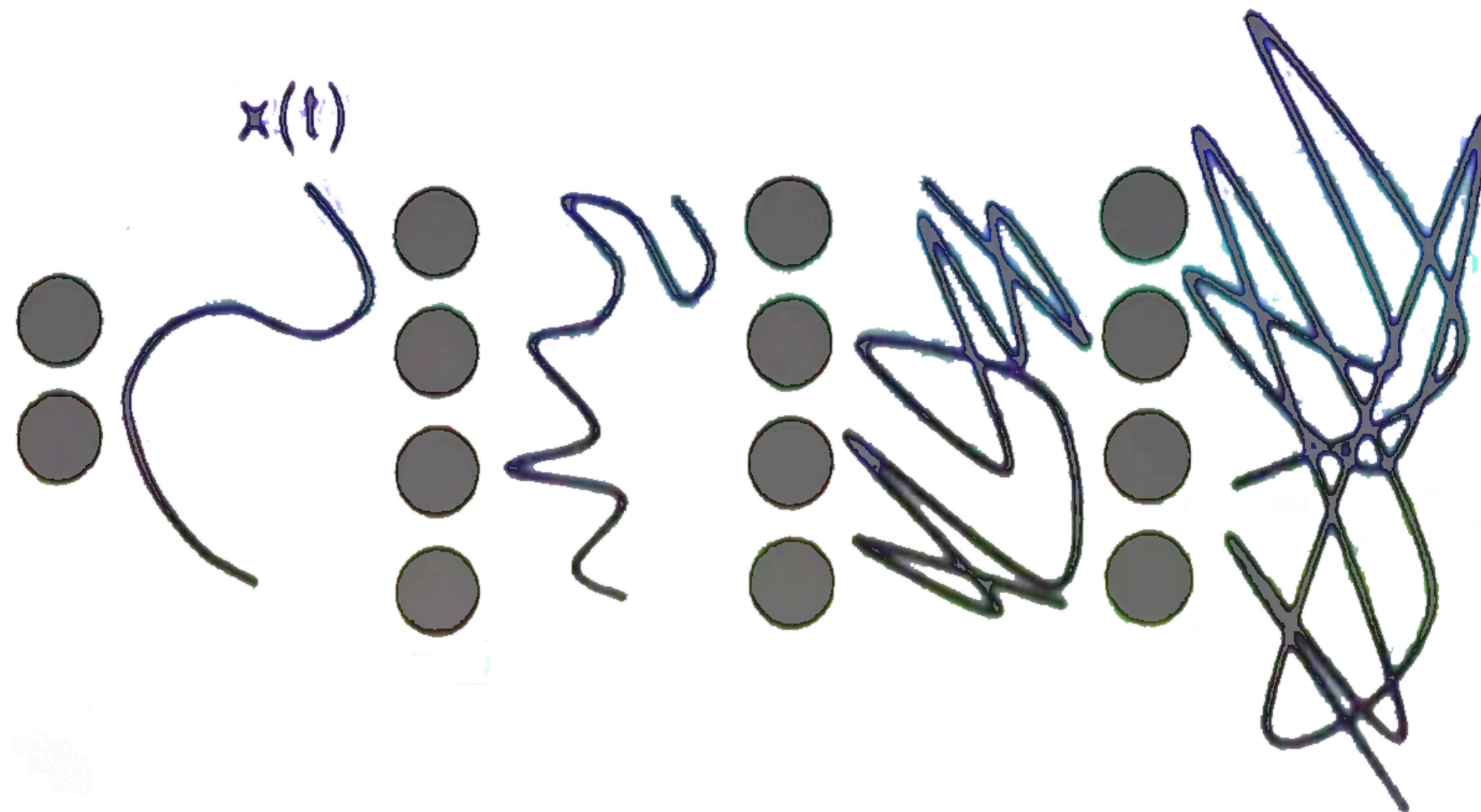
Rethink: What is Expressive Power?

- Consider an (one-dimensional) input trajectory



# Trajectory Length

Metric: Trajectory Length



How does the trajectory length increase?

# Trajectory Length

## What determines Trajectory Length?

For a network with depth  $d$ , width  $k$ , weights  $\sim \mathcal{N}(0, \sigma_w^2/k)$ , bias  $\sim \mathcal{N}(0, \sigma_b^2)$ , we have

(a)

$$\mathbb{E} \left[ l(z^{(d)}(t)) \right] \geq O \left( \frac{\sigma_w \sqrt{k}}{\sqrt{k+1}} \right)^d l(x(t))$$

for ReLUs



(b)

$$\mathbb{E} \left[ l(z^{(d)}(t)) \right] \geq O \left( \frac{\sigma_w \sqrt{k}}{\sqrt{\sigma_w^2 + \sigma_b^2 + k\sqrt{\sigma_w^2 + \sigma_b^2}}} \right)^d l(x(t))$$

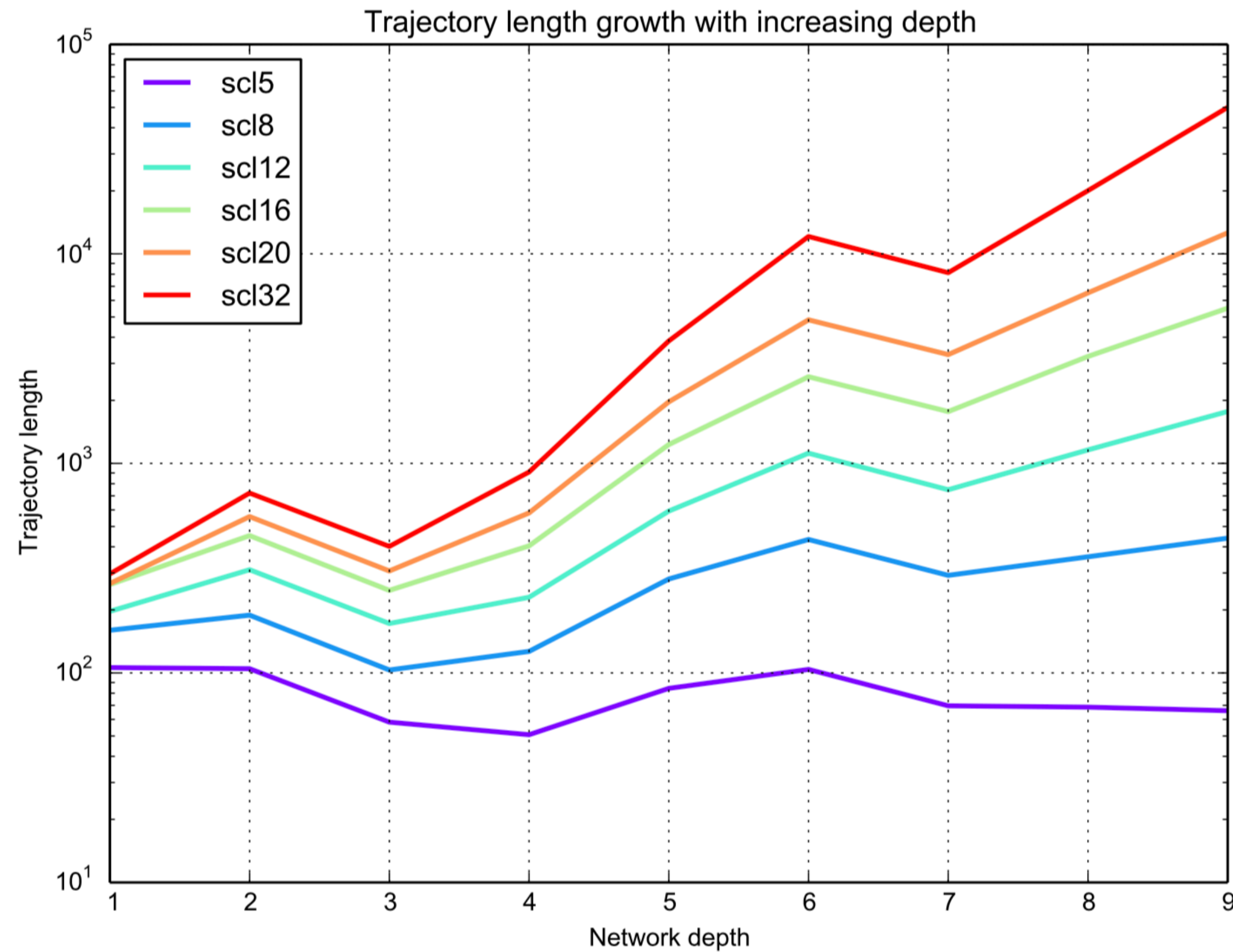
for hard tanh





# Trajectory Length

## Experiment Verification



Conv net on CIFAR-10, ReLU:

- Growth with depth exponentially
- Growth with  $\sigma_w$

# Trajectory Length

## Relationship with number of Linear Regions

For a hard tanh network with depth  $d$ ,  $n$  hidden layers, width  $k$ , weights  $\sim \mathcal{N}(0, \sigma_w^2/k)$ , bias  $\sim \mathcal{N}(0, \sigma_b^2)$ , we have

Review: Trajectory Length

(b)

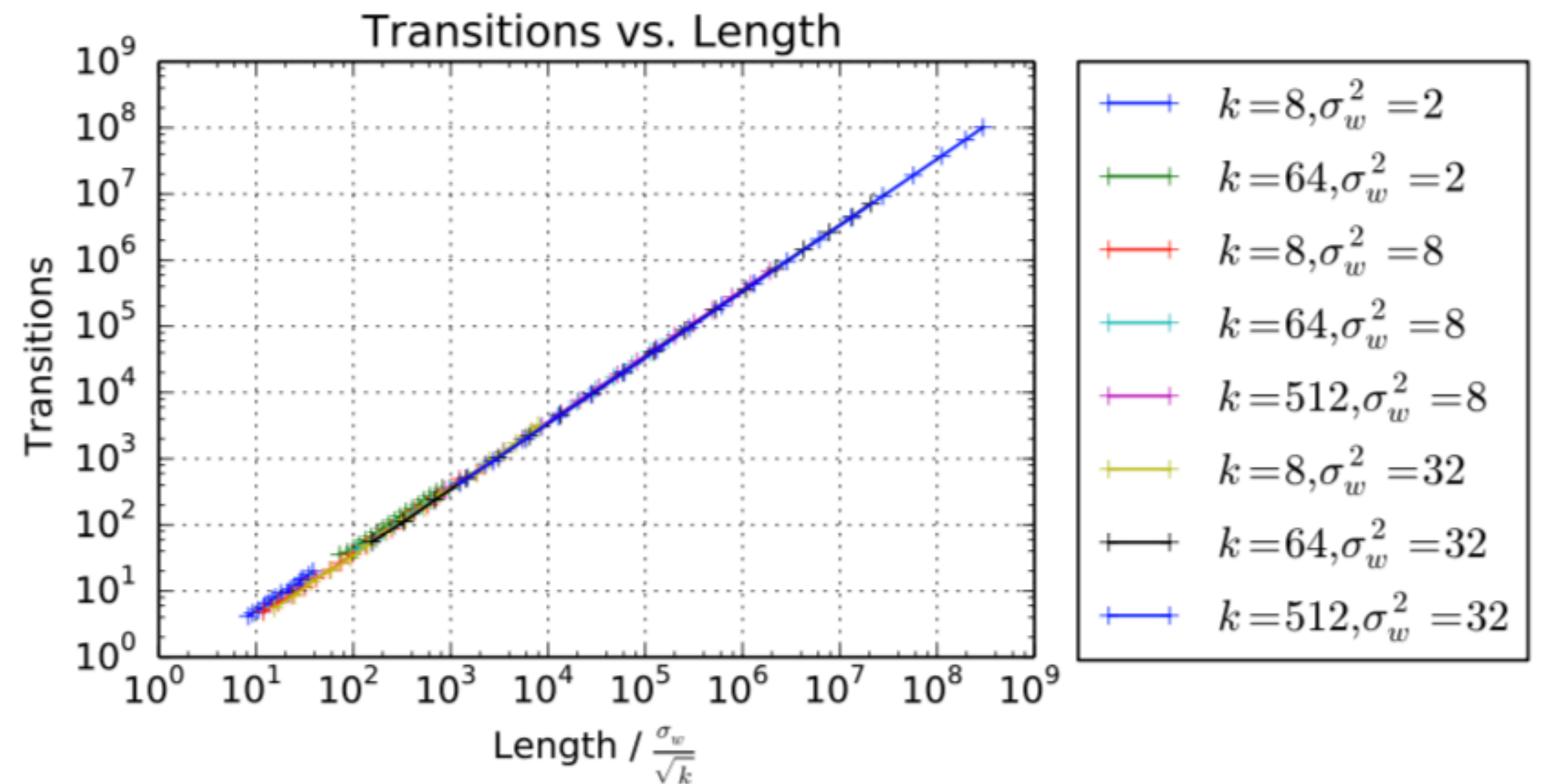
$$\mathbb{E} [l(z^{(d)}(t))] \geq O \left( \frac{\sigma_w \sqrt{k}}{\sqrt{\sigma_w^2 + \sigma_b^2 + k\sqrt{\sigma_w^2 + \sigma_b^2}}} \right)^d l(x(t))$$

for hard tanh

Transitions:

$$g(k, \sigma_w, \sigma_b, n) = O \left( \frac{\sqrt{k}}{\sqrt{1 + \frac{\sigma_b^2}{\sigma_w^2}}} \right)^n$$

Then  $\mathcal{T}(F_{A_{n,k}}(x(t); W) = O(g(k, \sigma_w, \sigma_b, n))$ .



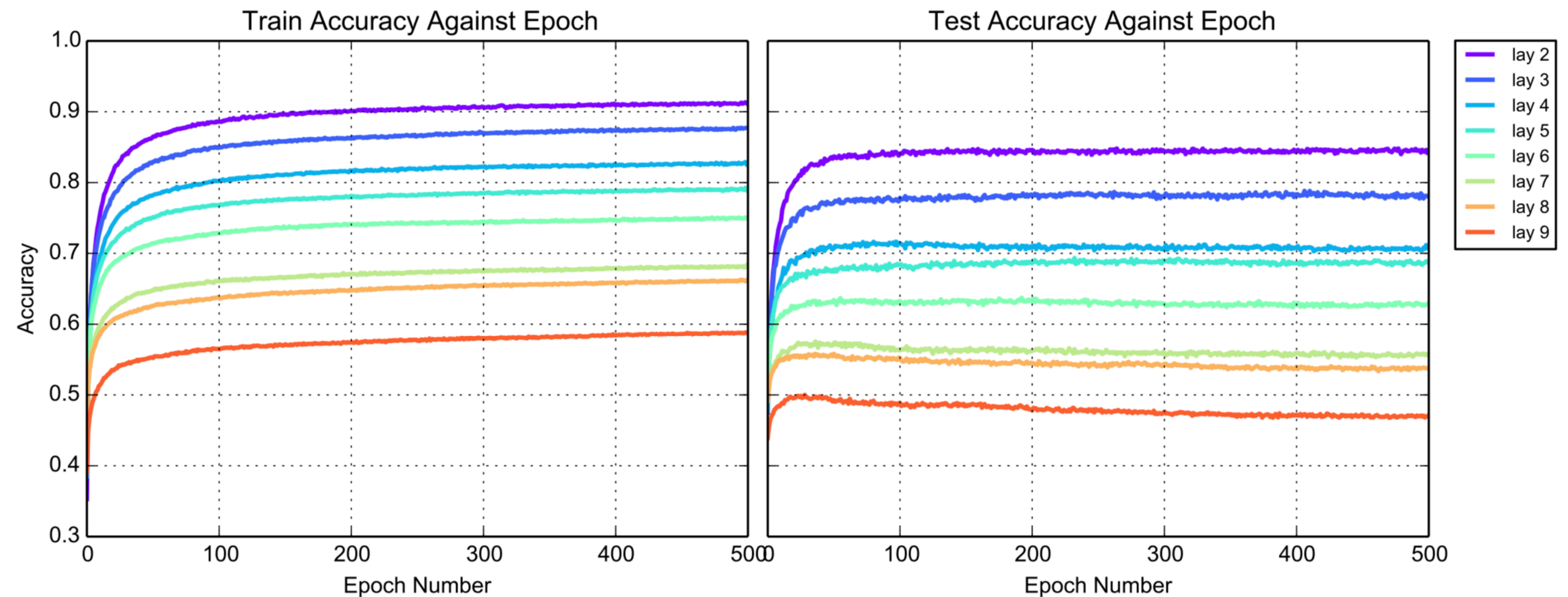
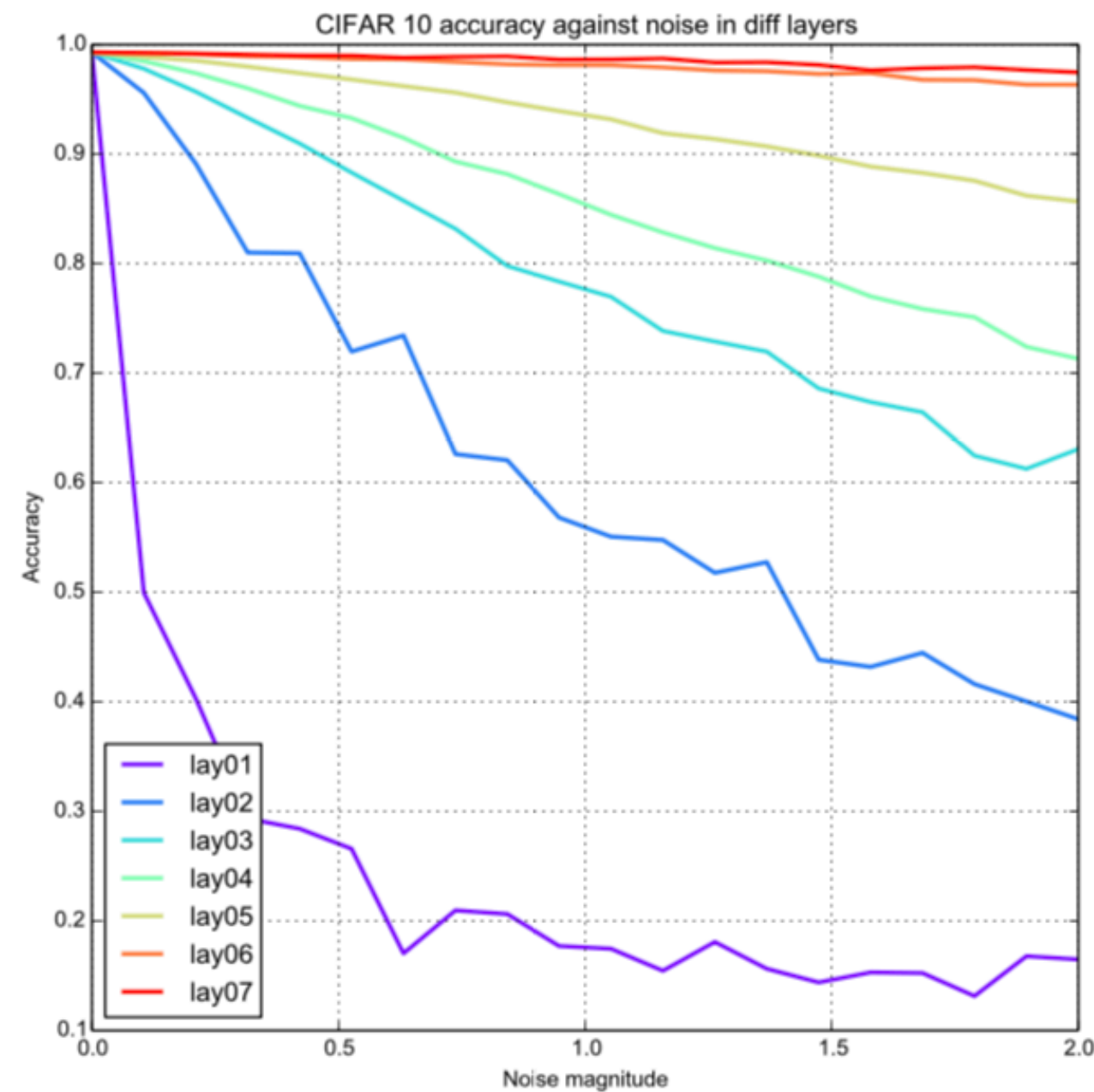
# Trajectory Length

## Insights: Trajectory and Stability

A perturbation at a layer grows exponentially in the *remaining depth* after that layer.

Add noise on different layers:

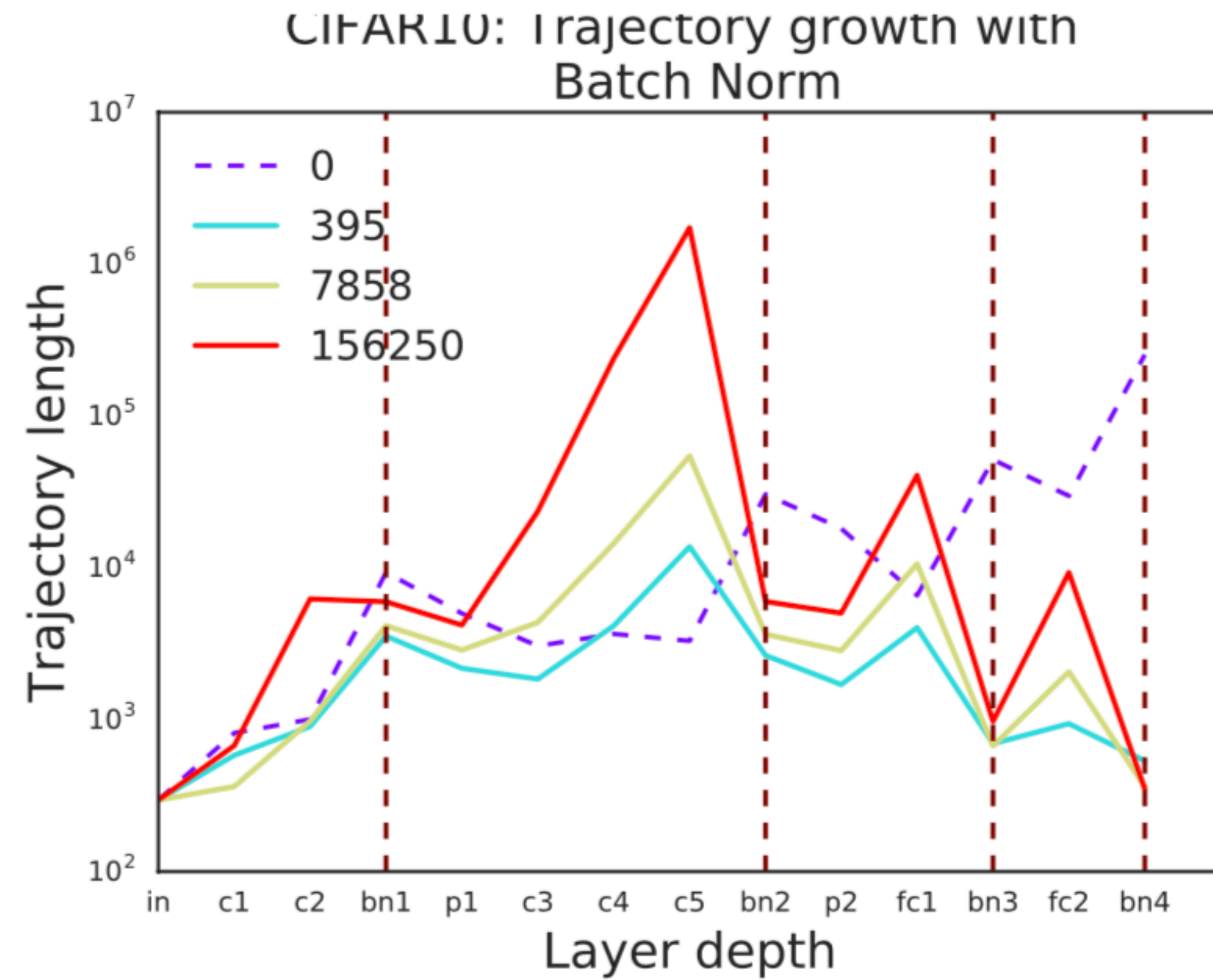
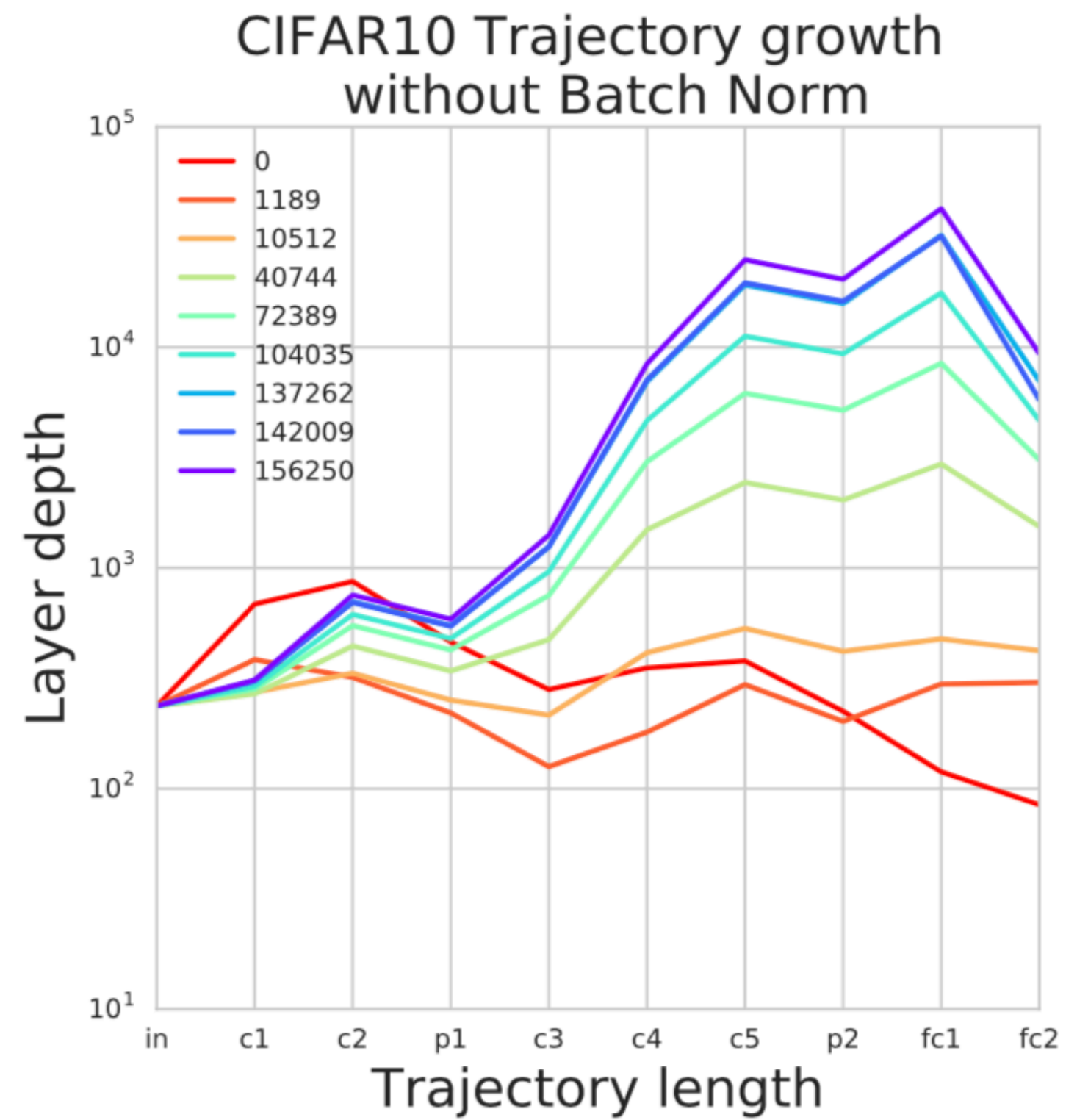
Only one layer trainable:





# Trajectory Length

## Insights: Trajectory and Batch Normalization



# Trajectory Length

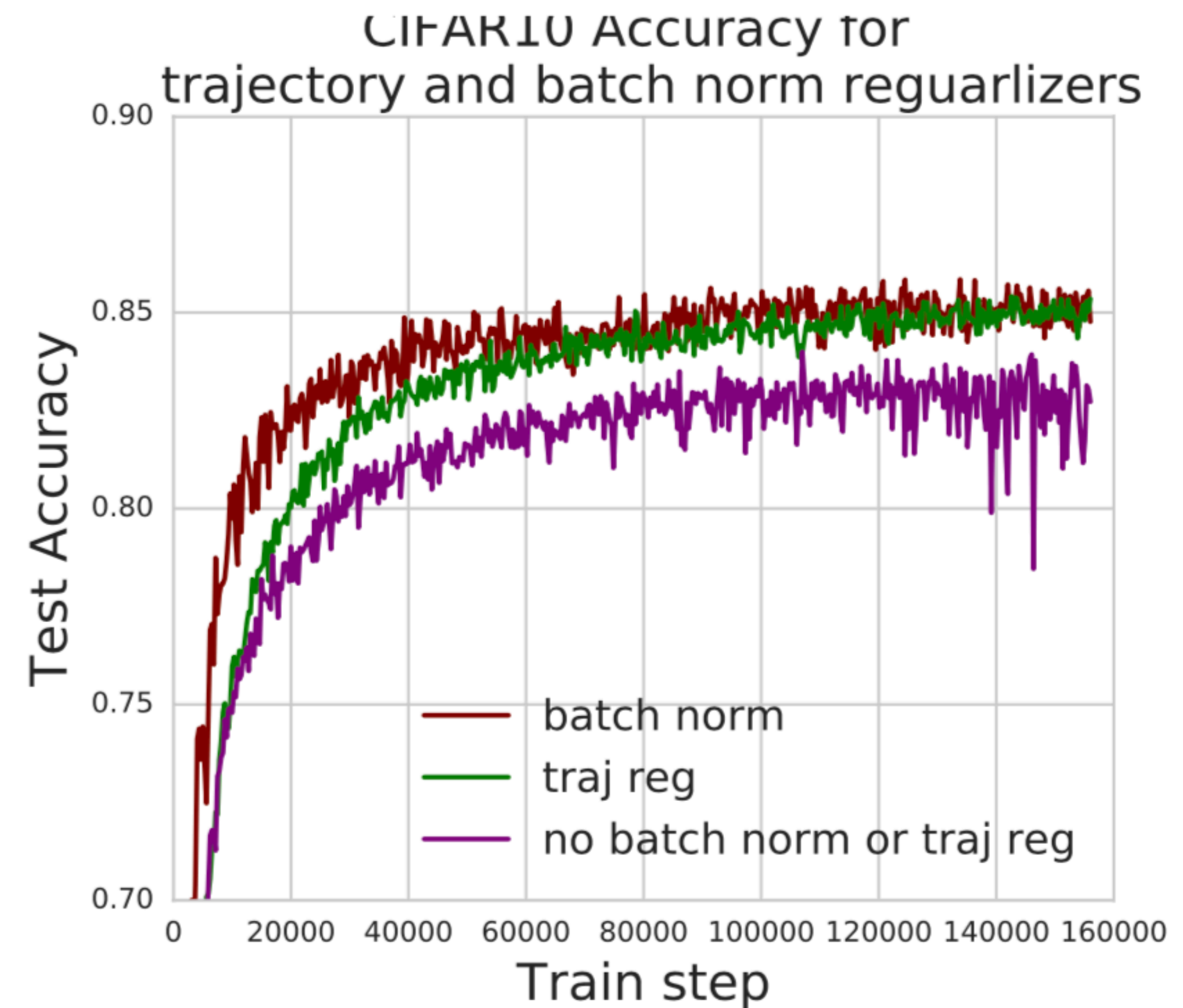
## Insights: Trajectory Normalization

Trajectory Normalization: Scale trajectory length directly.

Trajectory regularization layers: add to the loss

$$\lambda \frac{\text{currentlength}}{\text{originallength}}$$

In practice, compute the sum of distances between adjacent points in the mini-batch.



# Trajectory Length

## Review

1. What is it?
  2. How to measure?
  3. What determines it?
  4. Usage?
1. Mapping complexity
  2. Trajectory length
  3.  $\mathbb{E}[l(z^{(d)}(t))] \geq O\left(\frac{\sigma_w \sqrt{k}}{\sqrt{k+1}}\right)^d l(x(t))$  ,  $\mathbb{E}[l(z^{(d)}(t))] \geq O\left(\frac{\sigma_w \sqrt{k}}{\sqrt{\sigma_w^2 + \sigma_b^2 + k\sqrt{\sigma_w^2 + \sigma_b^2}}}\right)^d l(x(t))$
  4. Stability, Batch-norm, Traj-reg

## Question?

# Contents



**Introduction**

---



**Activation Pattern**

---



**Trajectory Length**

---



**Conclusion**

---

# Conclusion

## Related Materials

- Paper and Supplementary Link: <http://proceedings.mlr.press/v70/raghu17a.html>
- Presentation Video: <https://vimeo.com/237276052>
- ICLR 2017 discussion: <https://openreview.net/forum?id=B1TTpYKgx>

# Conclusion

## Further Work

- What about other activation functions?
  - Their previous paper [1] talked about it. It combines Riemannian geometry with the mean field theory of high dimensional chaos to study it. (What are they? :-)
- What about setting the input as a plane or other hyper-space, instead of only trajectory?
- What if the network is not regular (i.e. the width is not the same?)
- Is there any other proper metric for expressive power?

[1] B. Poole, S. Lahiri, M. Raghu, J. Sohl-Dickstein and S. Ganguli. Exponential expressivity in deep neural networks through transient chaos. In NeurIPS, 2016

# Conclusion

Thank you!

Any questions or evaluations?