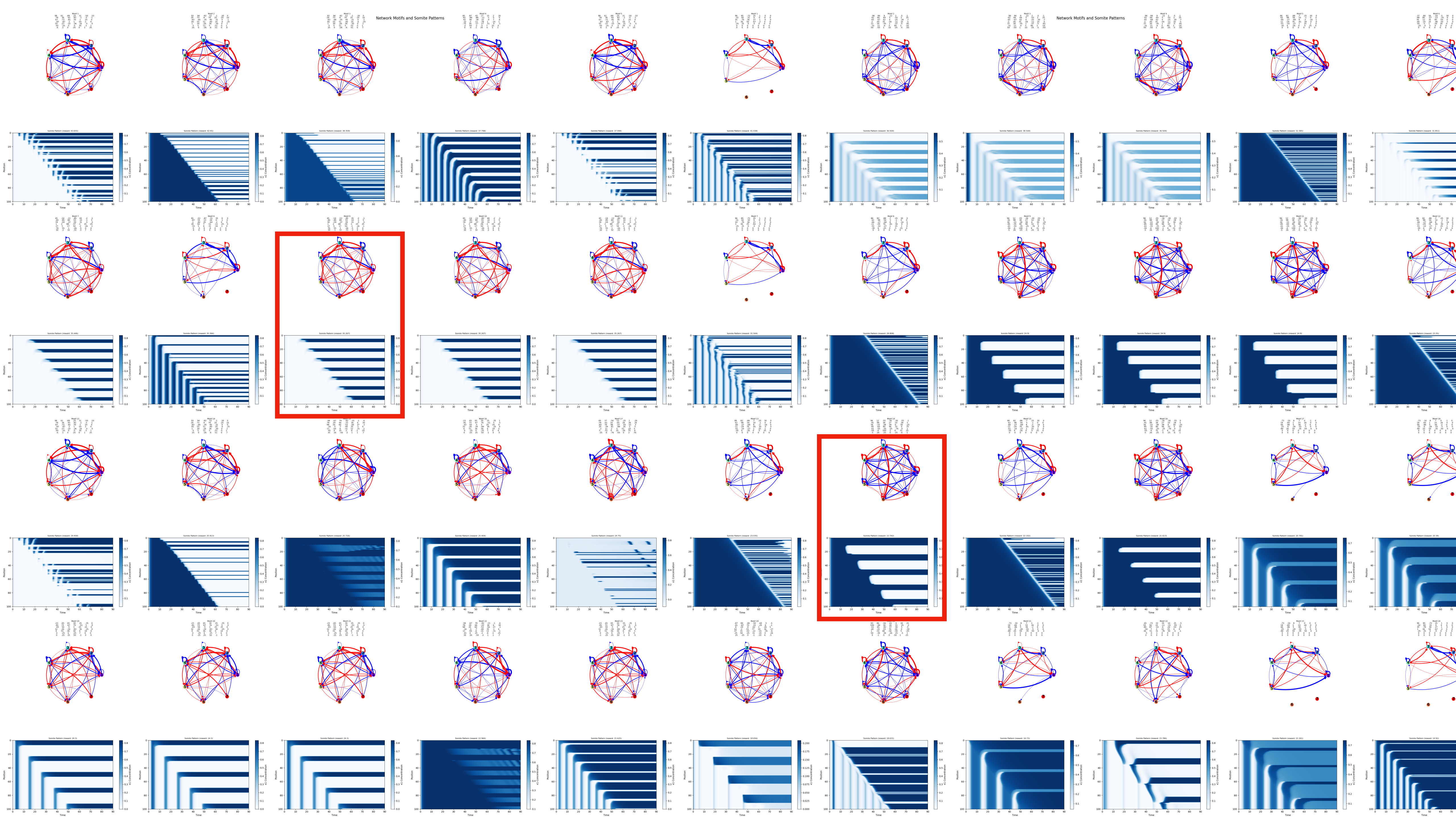


How to prune the network ?

THE LOTTERY TICKET HYPOTHESIS:
FINDING SPARSE, TRAINABLE NEURAL NETWORKS

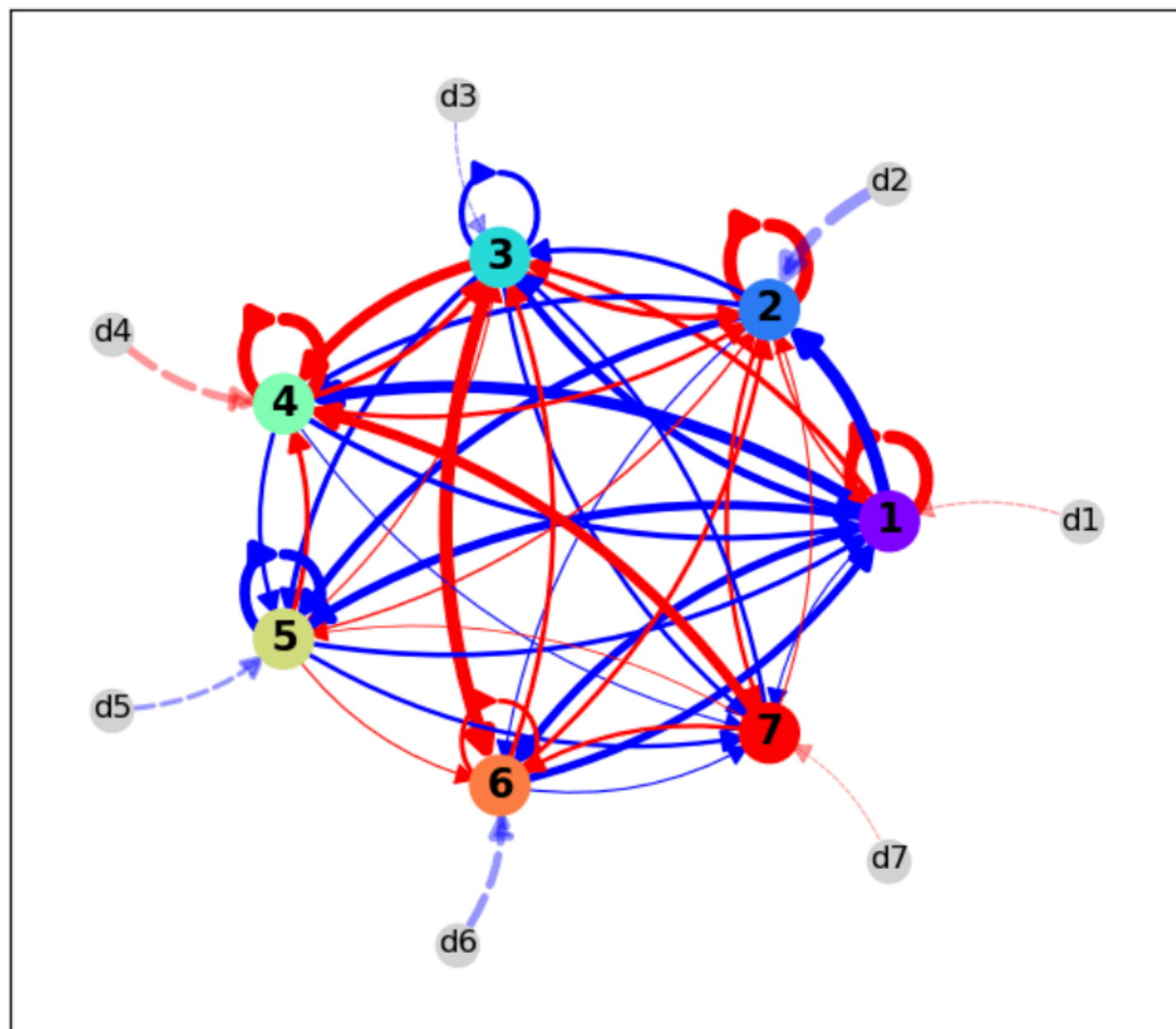
Jonathan Frankle
MIT CSAIL
jfrankle@csail.mit.edu

Michael Carbin
MIT CSAIL
mcarbin@csail.mit.edu



Motif 9						
96	12	-55	-36	-30	-55	0
-91	98	38	27	10	28	5
36	-33	-38	38	6	30	-25
-85	-35	80	96	32	0	75
-61	-55	-35	-25	-76	0	1
-56	-6	100	-1	6	31	25
-5	26	-25	-5	-25	-6	-1

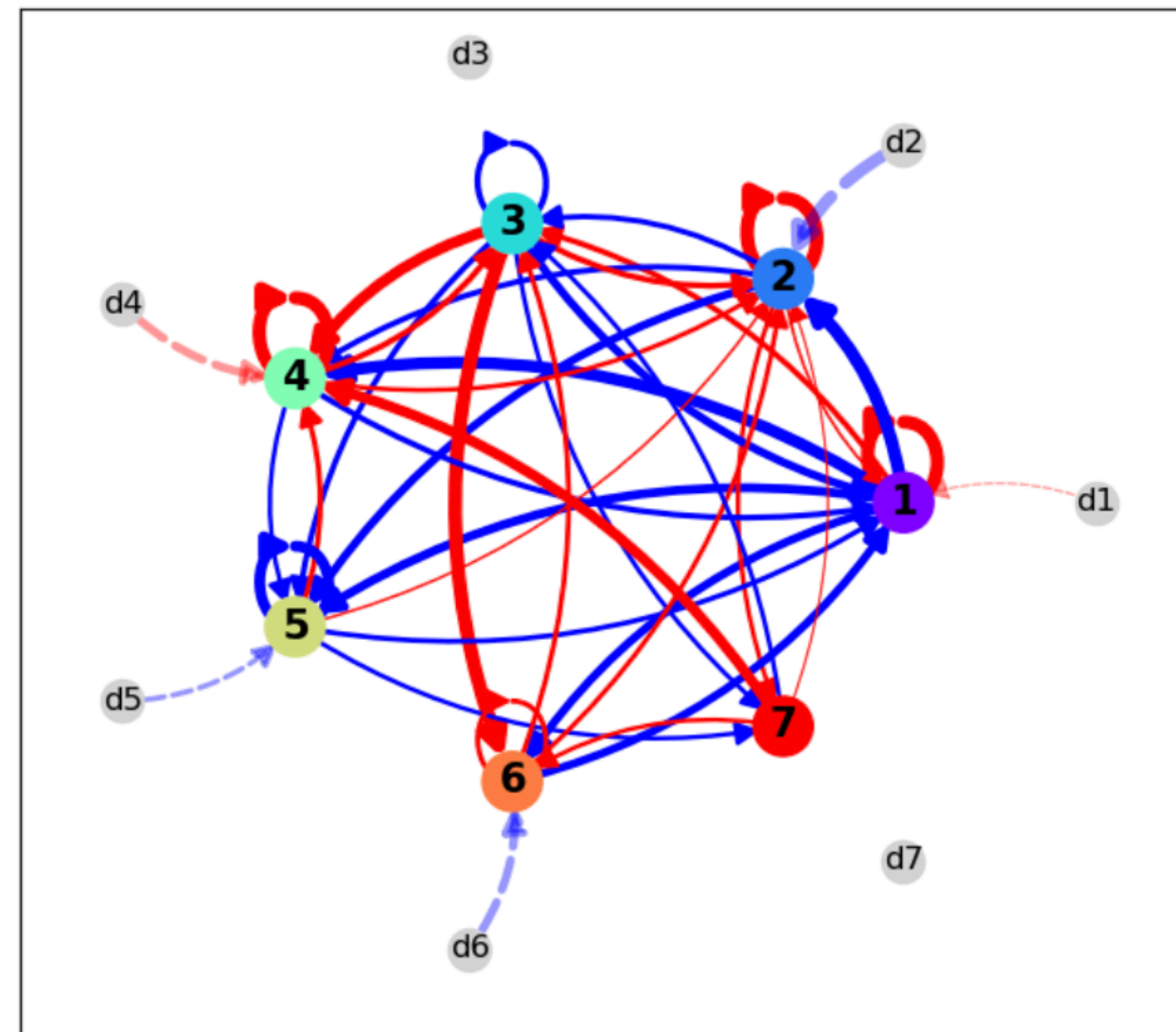
7-Node Network Motif



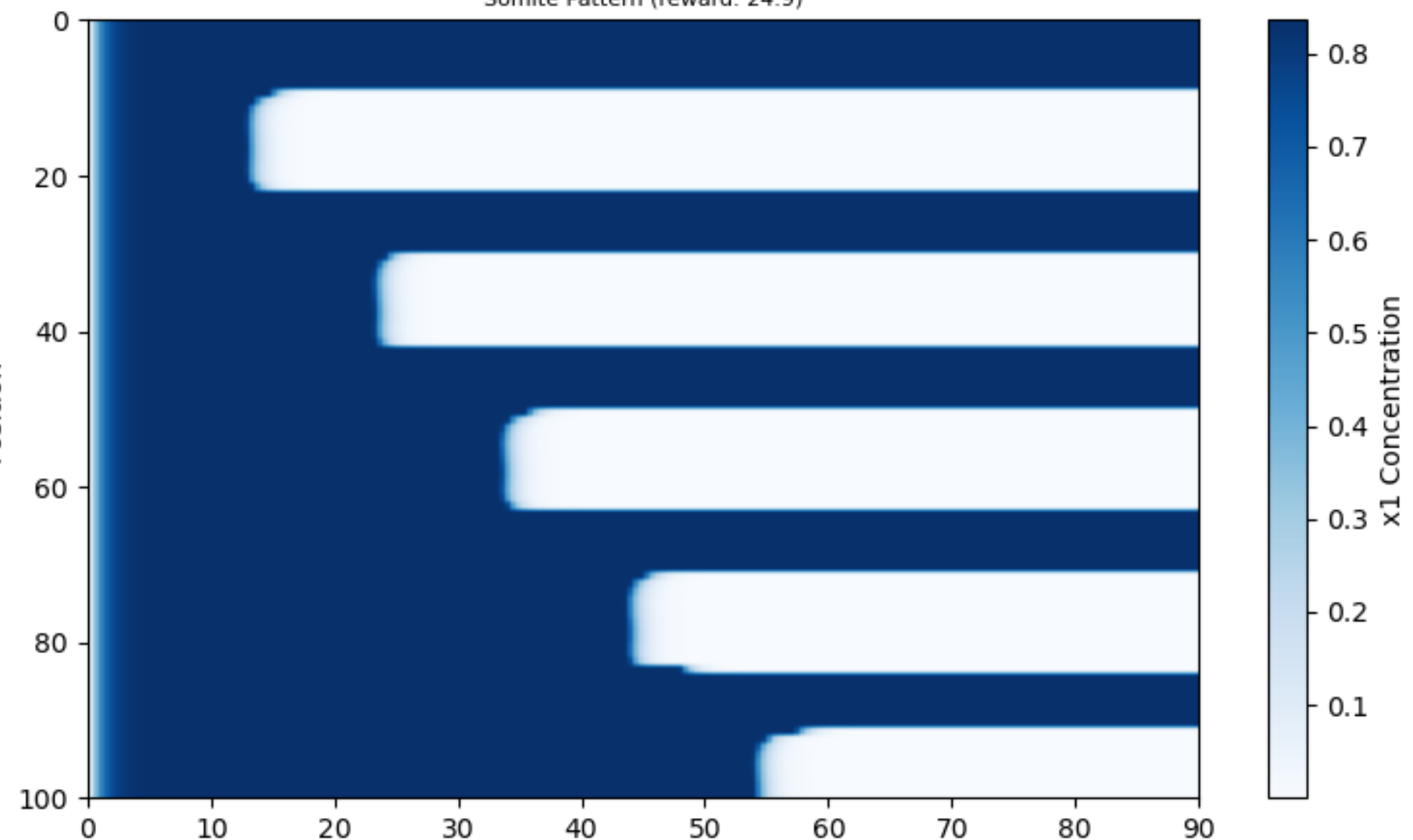
prunning



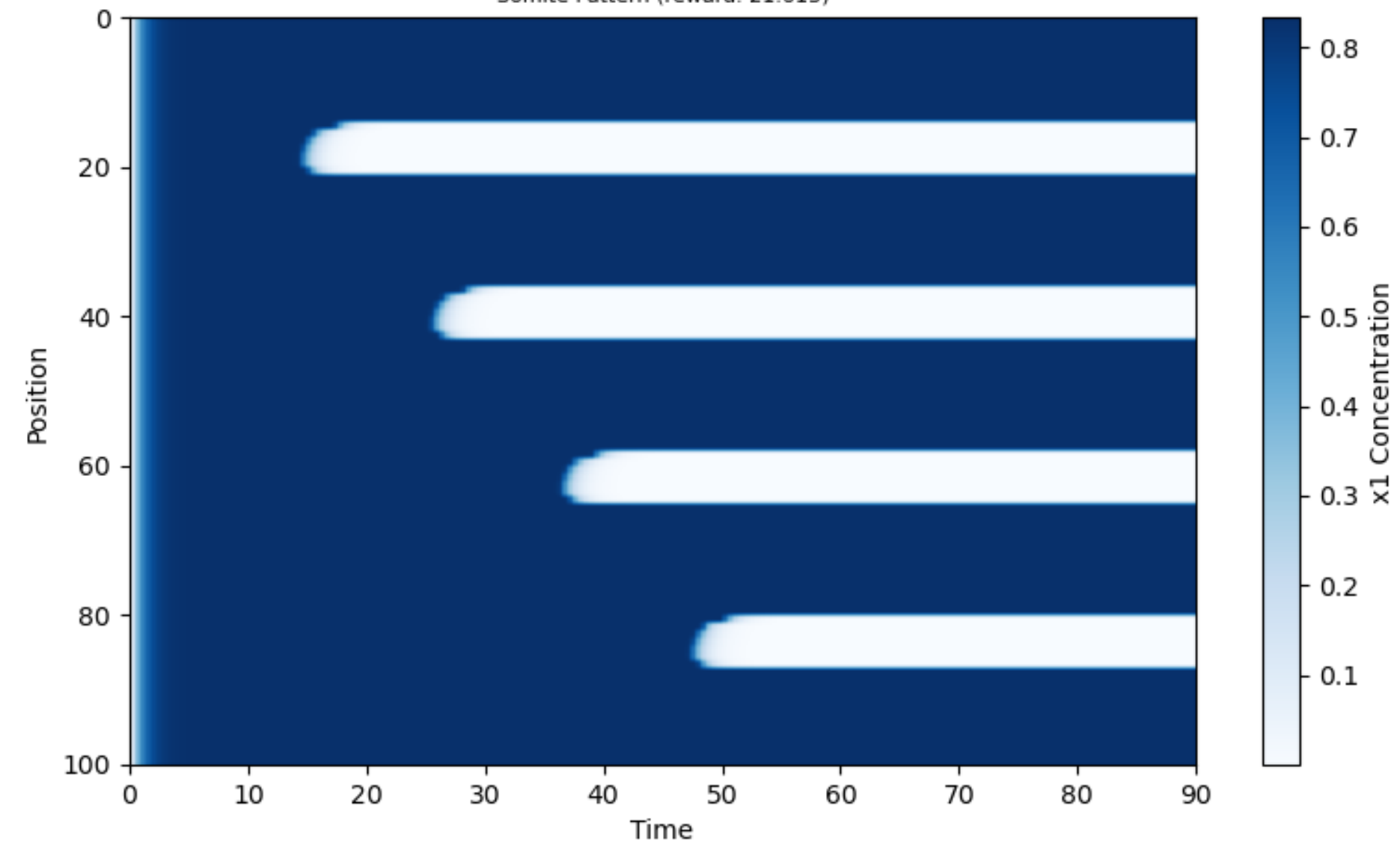
7-Node Network Motif



Somite Pattern (reward: 24.9)



Somite Pattern (reward: 21.613)

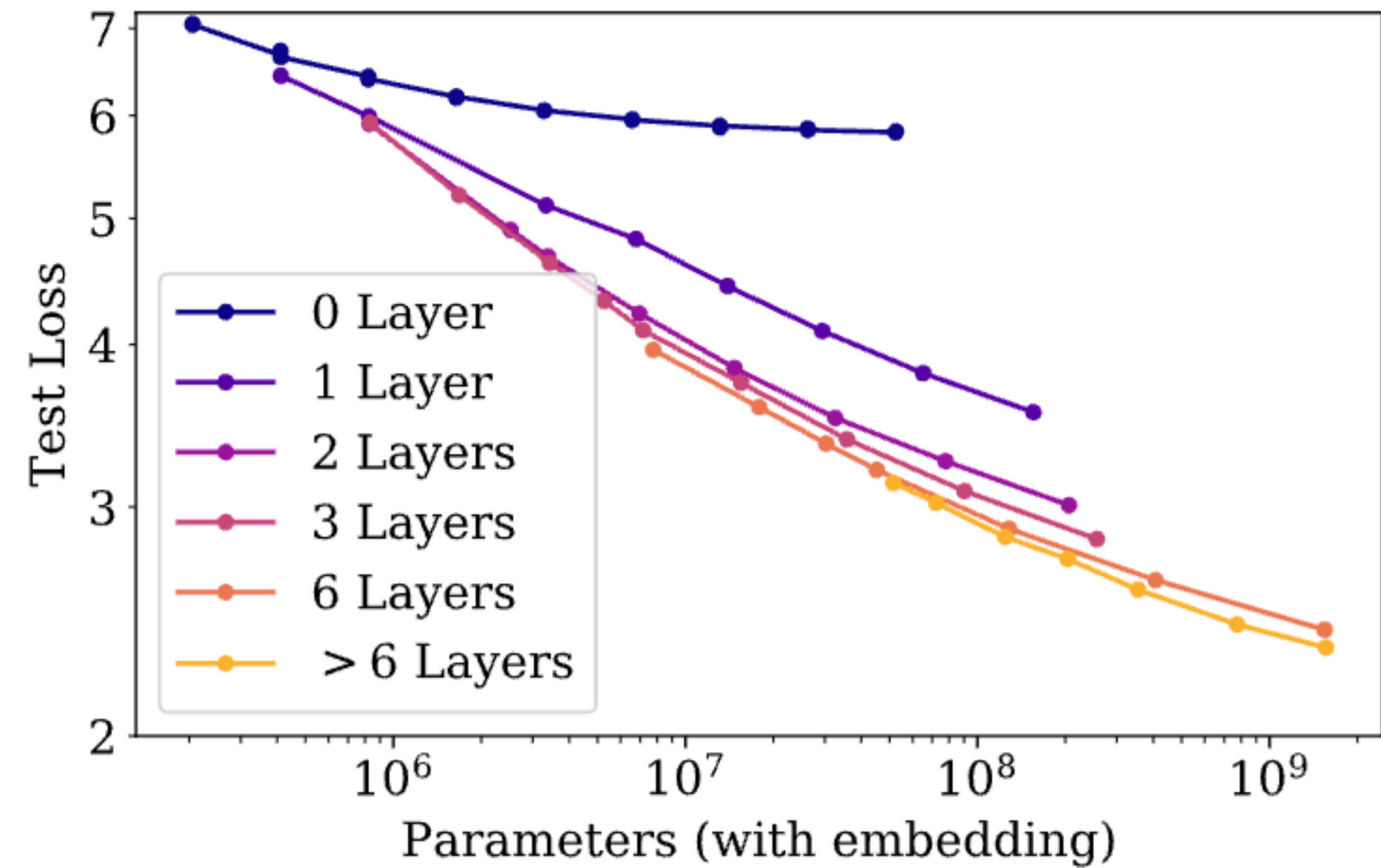
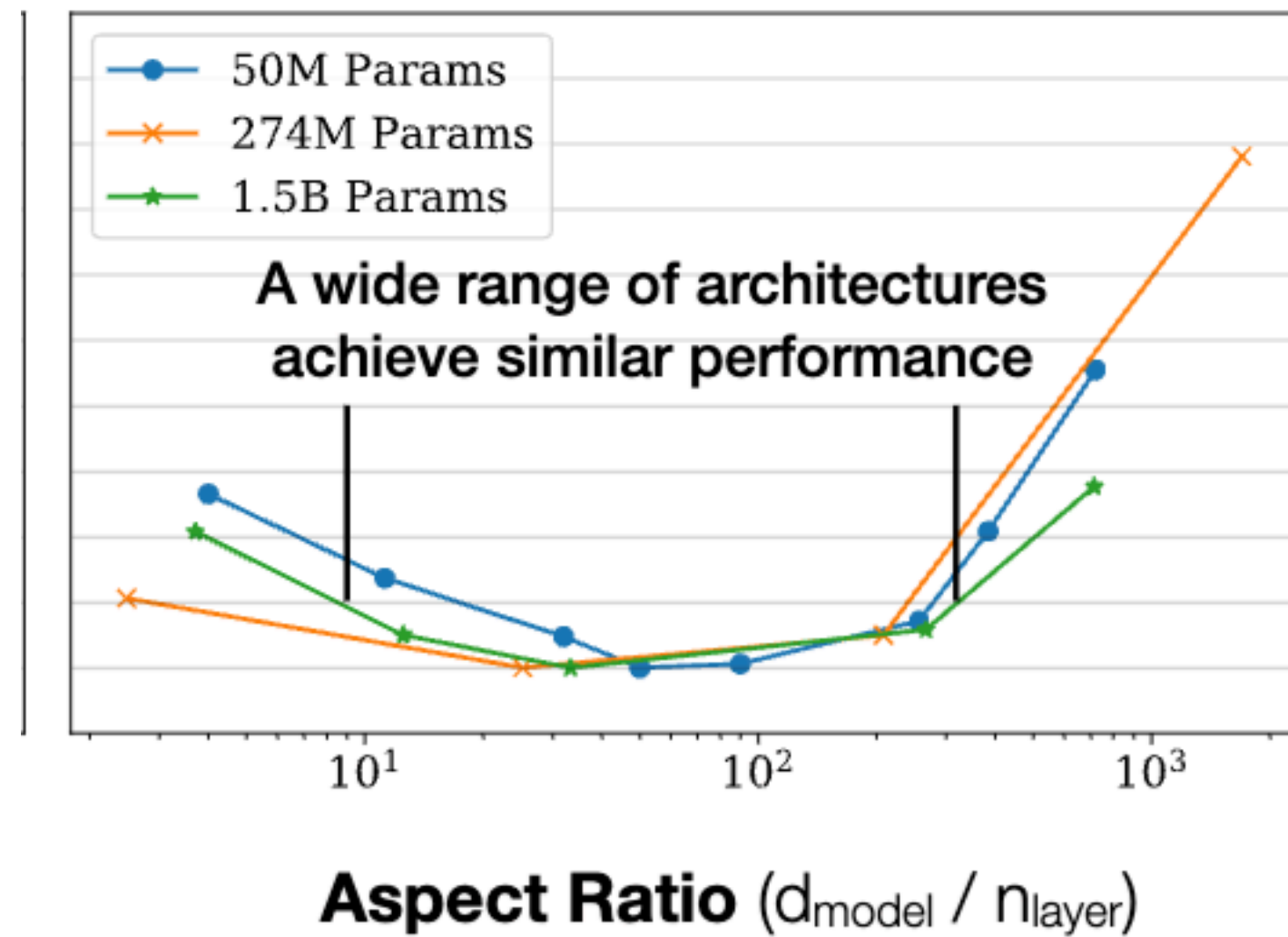


How to prune the network ?

Why prune the network ? / Why not to prune the network

Language Models are Few-Shot Learners

Scaling Laws for Neural Language Models



Learning both Weights and Connections for Efficient Neural Networks

Song Han
Stanford University
songhan@stanford.edu

Jeff Pool
NVIDIA
jpool@nvidia.com

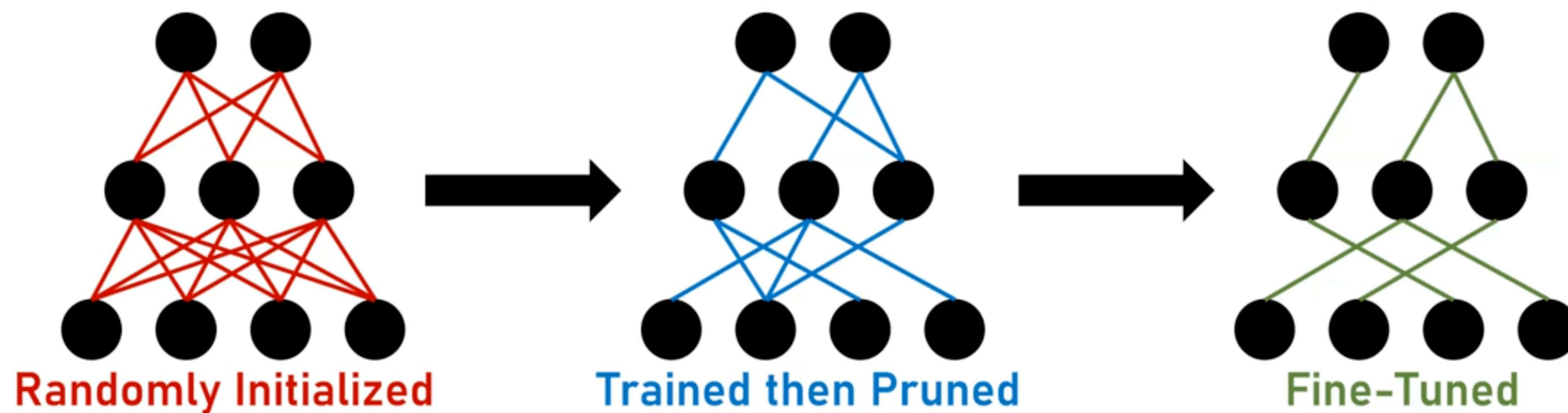
John Tran
NVIDIA
johntran@nvidia.com

William J. Dally
Stanford University
NVIDIA
dally@stanford.edu

2015

Pruning **weights** with lowest magnitude

- 1) Train the network
- 2) Remove superfluous structure
- 3) Fine-tune the network
- 4) Optionally: prune and fine-tune iteratively



Accuracy dropped

THE LOTTERY TICKET HYPOTHESIS:
FINDING SPARSE, TRAINABLE NEURAL NETWORKS

Jonathan Frankle
MIT CSAIL
jfrankle@csail.mit.edu

Michael Carbin
MIT CSAIL
mcarbin@csail.mit.edu



The Lottery Ticket Hypothesis :

The neural networks we typically train have **subnetworks** (at non-trivial sparsities) at **initialization** that can train to full accuracy in the same number of steps as the original network.

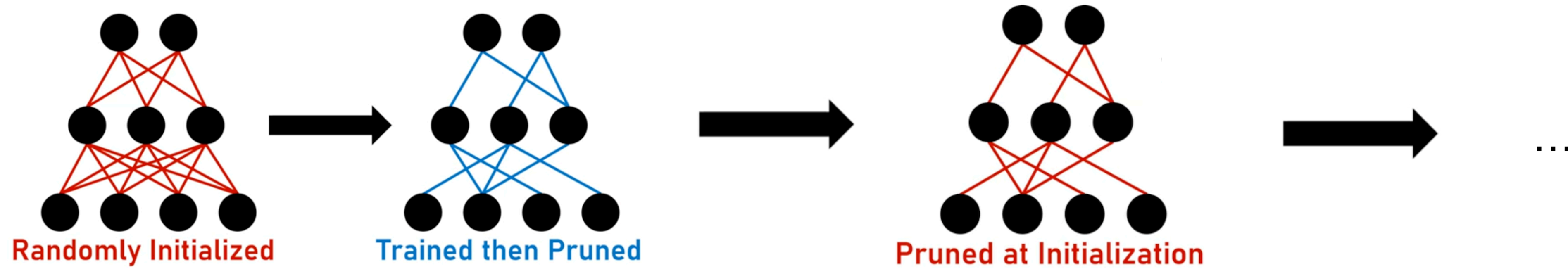
THE LOTTERY TICKET HYPOTHESIS: FINDING SPARSE, TRAINABLE NEURAL NETWORKS

Jonathan Frankle
MIT CSAIL
jfrankle@csail.mit.edu

Michael Carbin
MIT CSAIL
mcarbin@csail.mit.edu

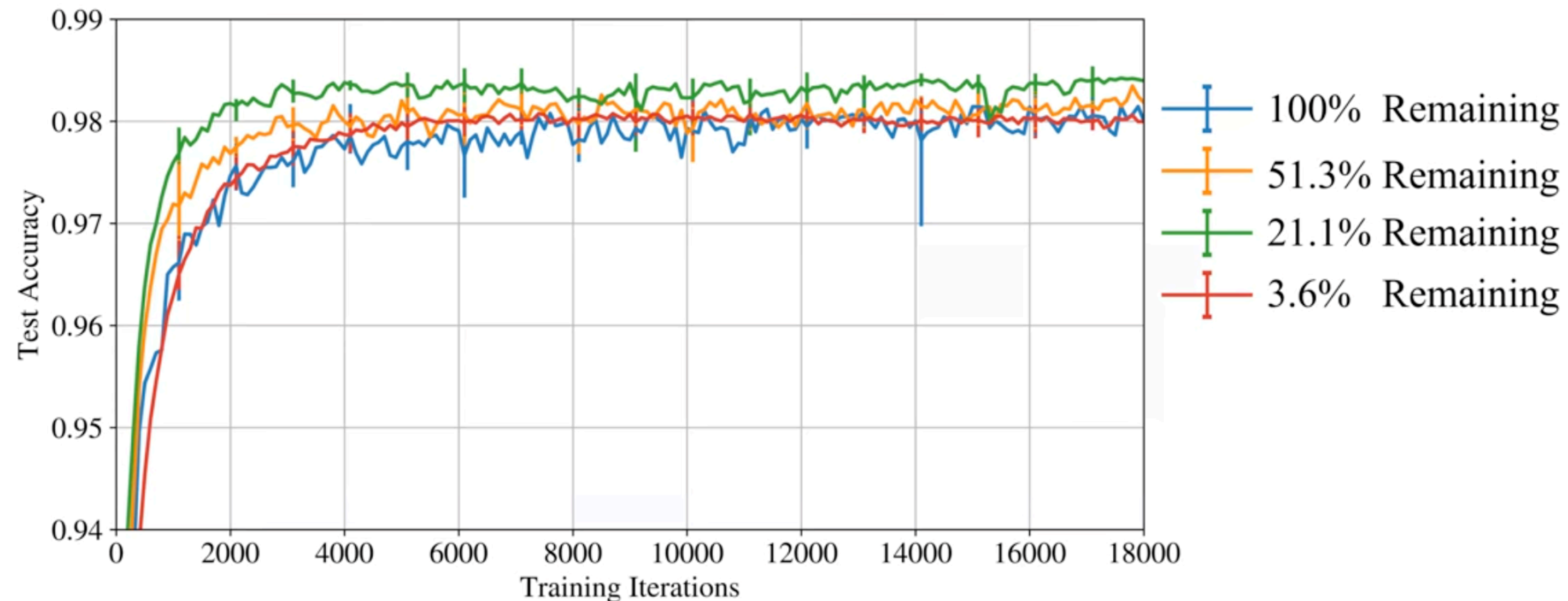
The Importance of Initialization

- 1) Train the network
- 2) Remove superfluous structure
- 3) Reset each weight to its original initialization
- 4) Train it to convergence

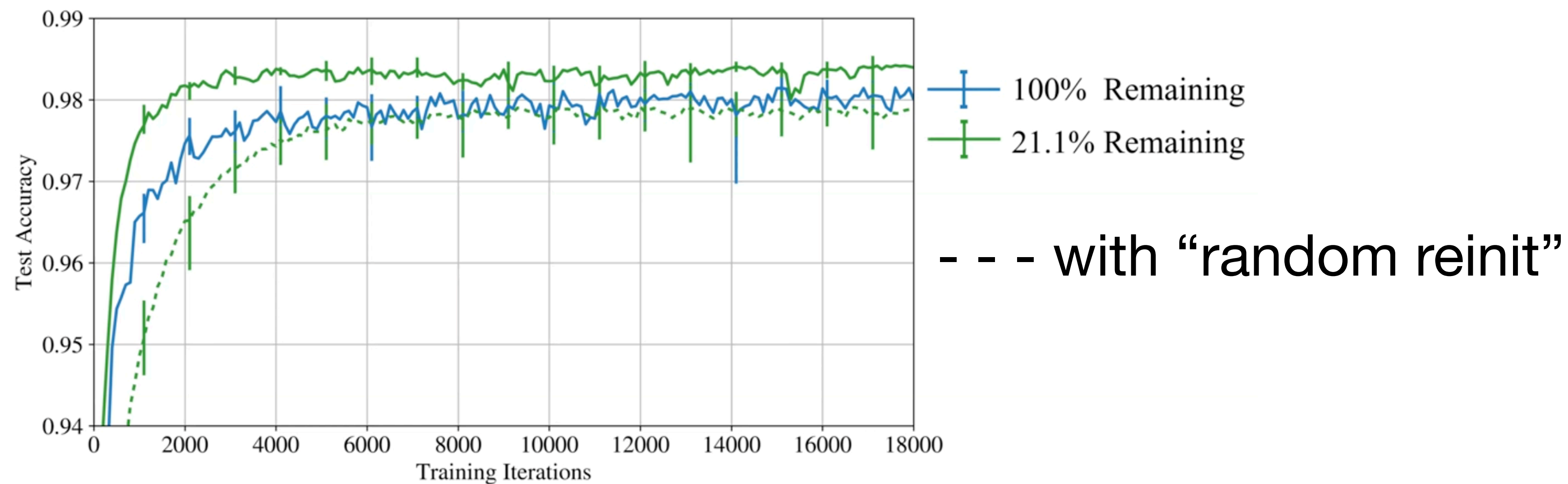


Reset weights to its original initialization

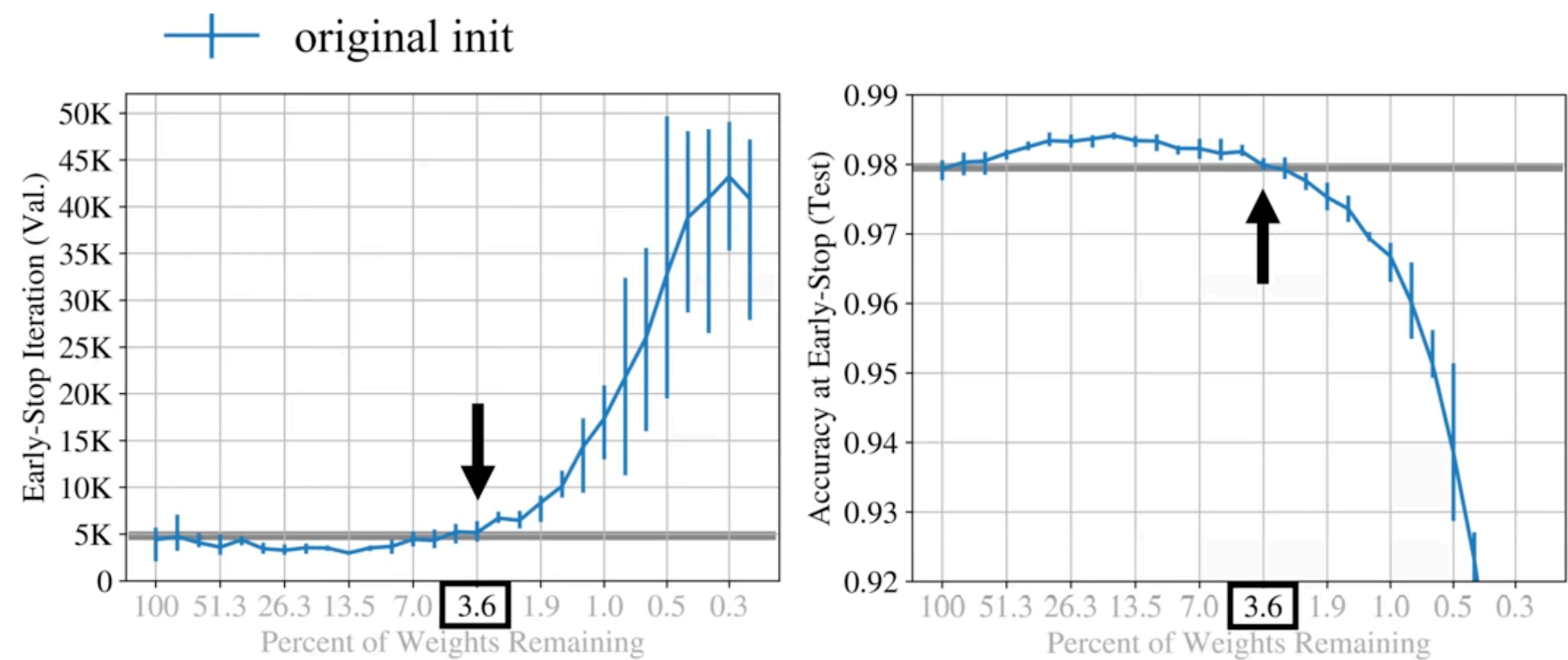
Retrain



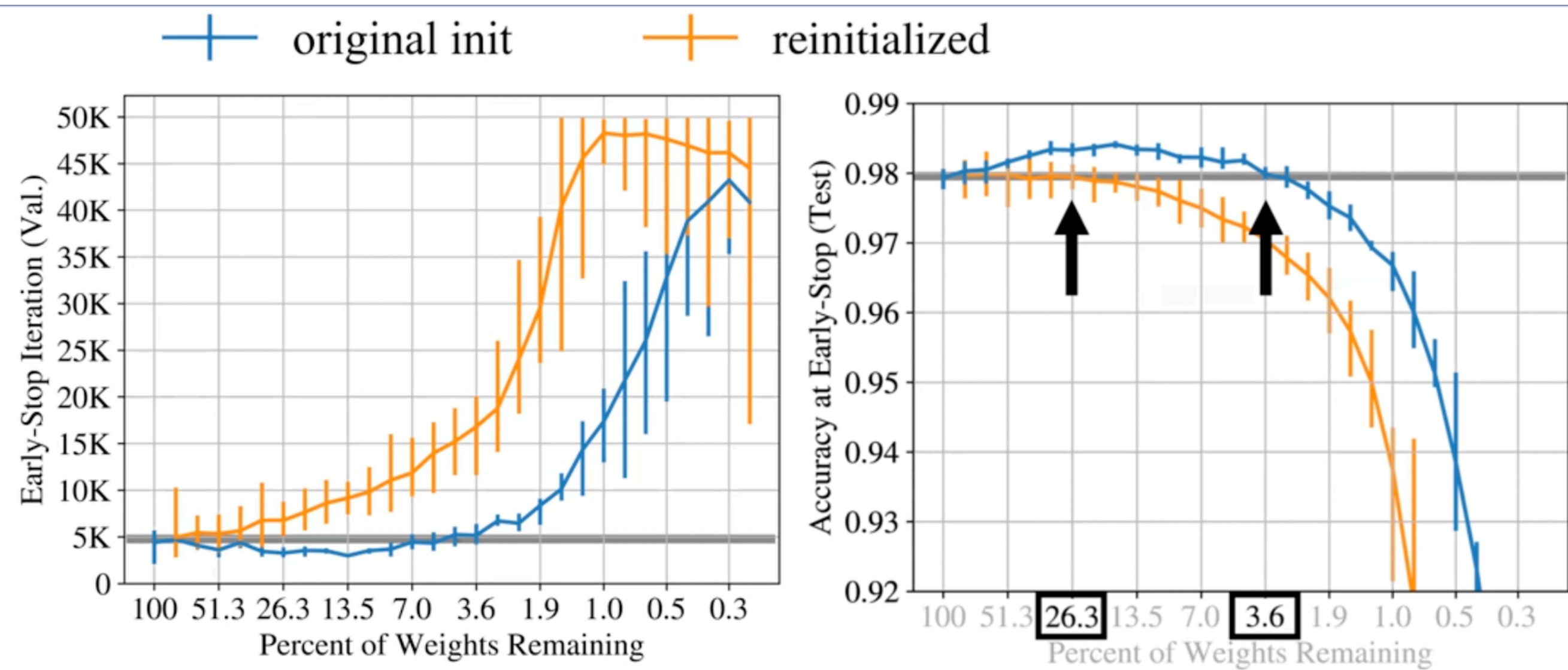
LeNet 300-100-10 for MNIST fully-connected 300K parameters



LeNet 300-100-10 for MNIST fully-connected 300K parameters

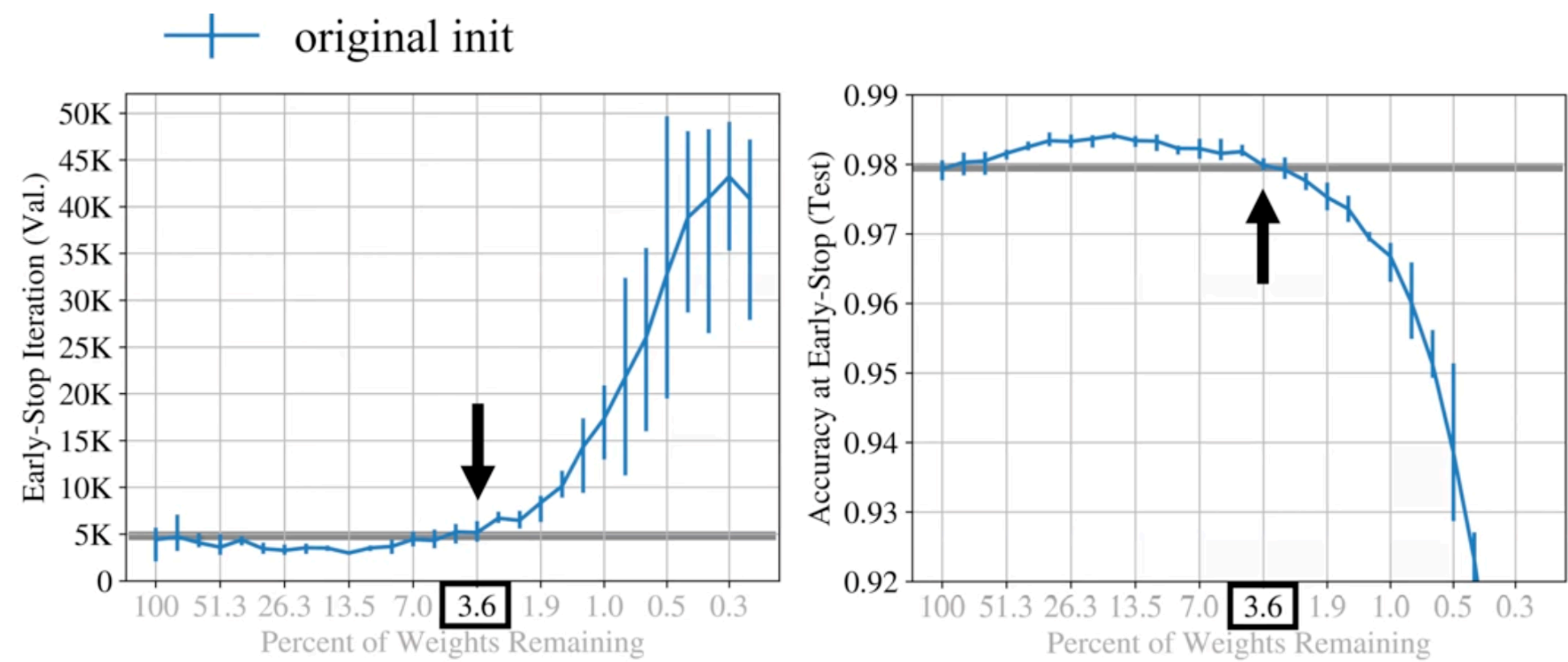


LeNet 300-100-10 for MNIST fully-connected 300K parameters

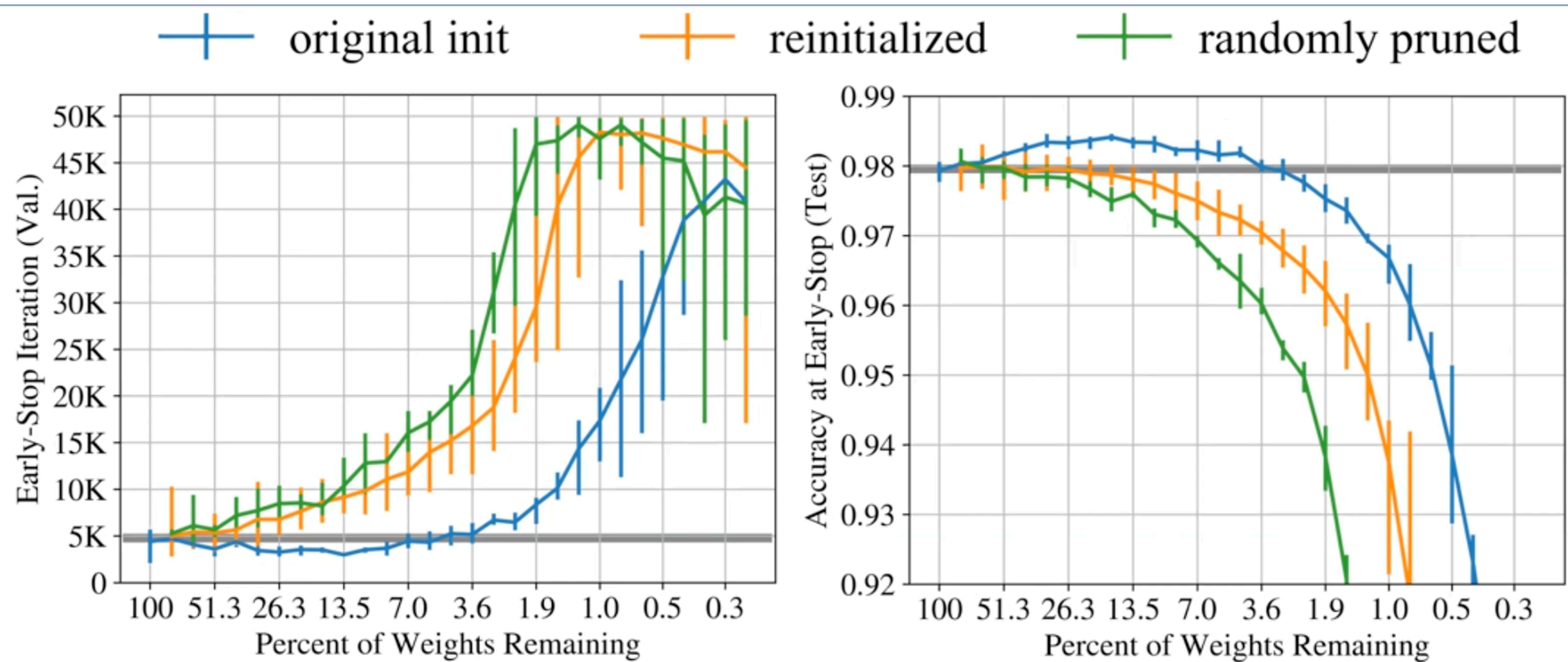


LeNet 300-100-10 for MNIST fully-connected 300K parameters

initialization matters



LeNet 300-100-10 for MNIST fully-connected 300K parameters



LeNet 300-100-10 for MNIST fully-connected 300K parameters

subnetwork
architecture matters

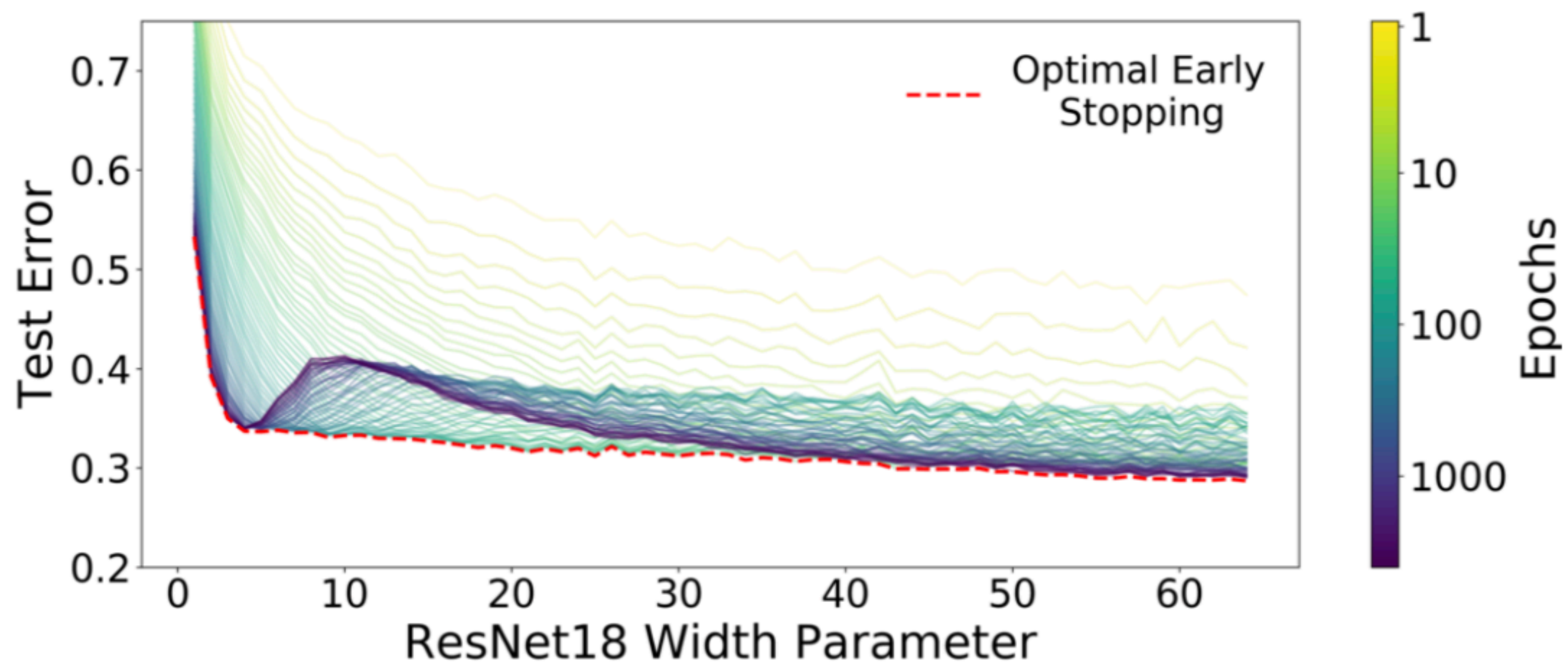


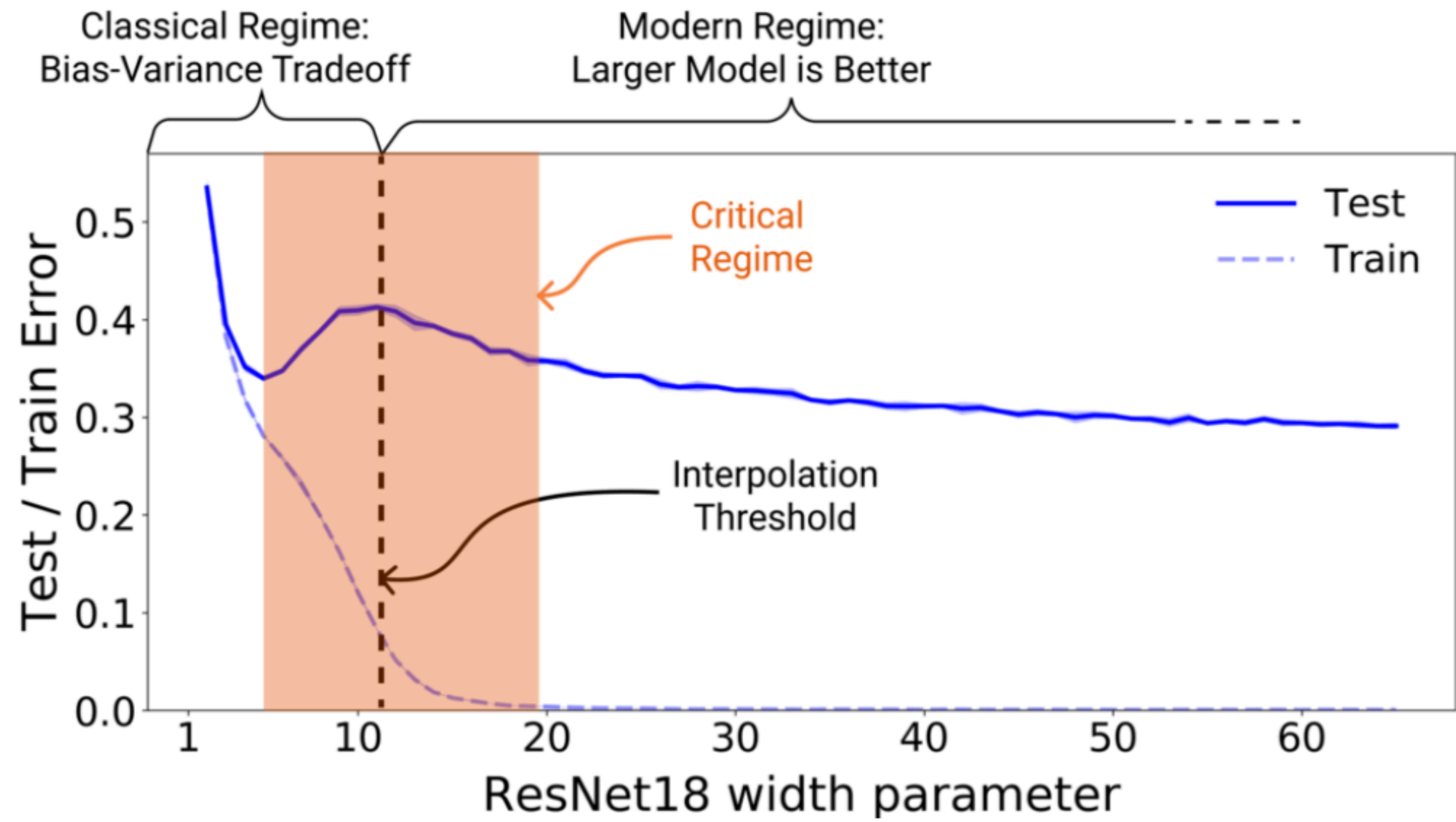
The Lottery Ticket Hypothesis :

The neural networks we typically train have **subnetworks** (at non-trivial sparsities) at **initialization** that can train to full accuracy in the same number of steps as the original network.

**If we can use this smaller network,
then why do we even bother with the larger one ?**

double descent





statistical learning theory
 The simplest model really is
 the best generalizer.

Occam's razor, also known as the principle of parsimony, suggests that when faced with multiple explanations for a phenomenon, the simplest one is usually the most likely to be correct.



Ockham chooses a razor

THE LOTTERY TICKET HYPOTHESIS: FINDING SPARSE, TRAINABLE NEURAL NETWORKS

Jonathan Frankle

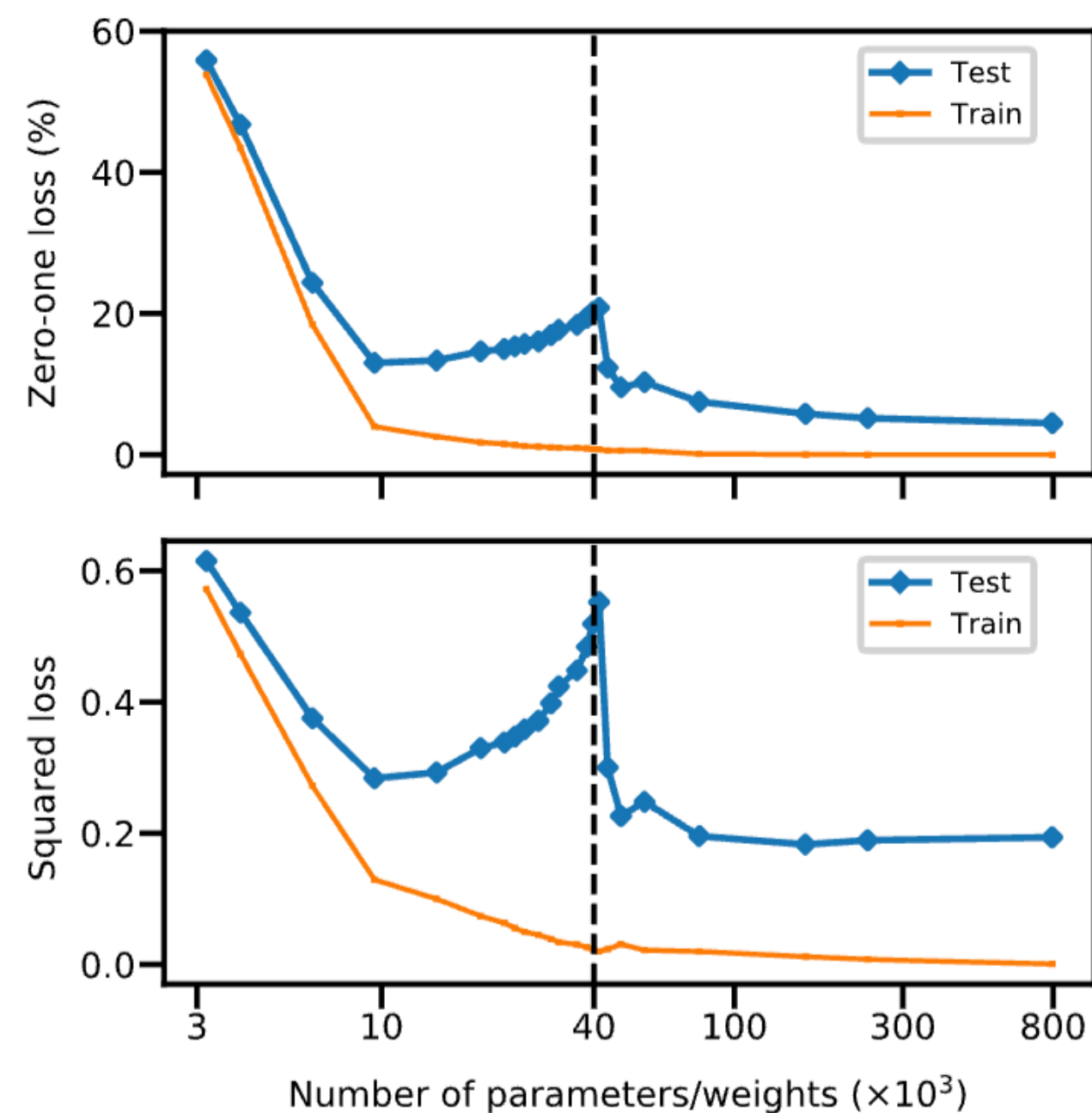
MIT CSAIL

jfrankle@csail.mit.edu

Michael Carbin

MIT CSAIL

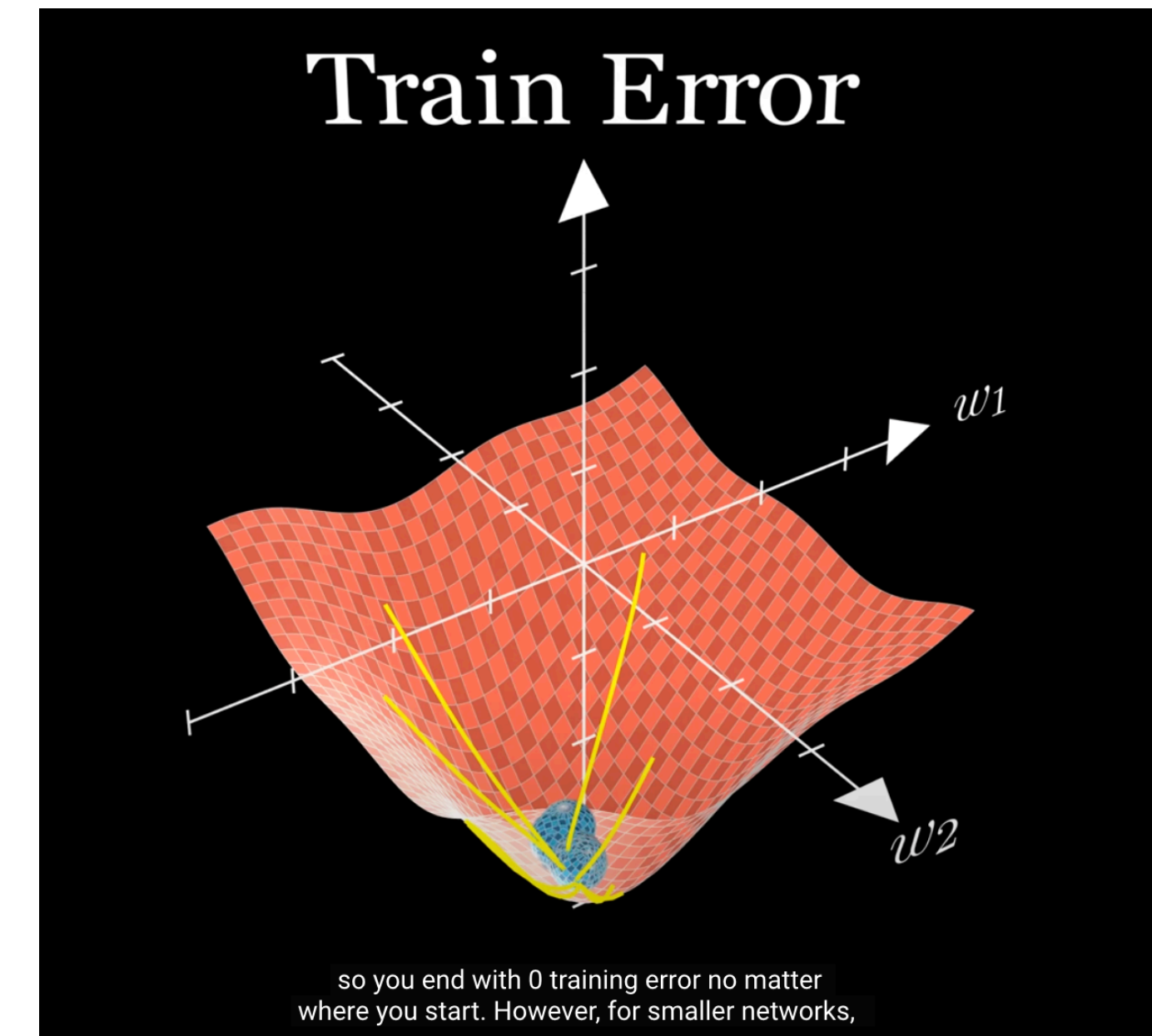
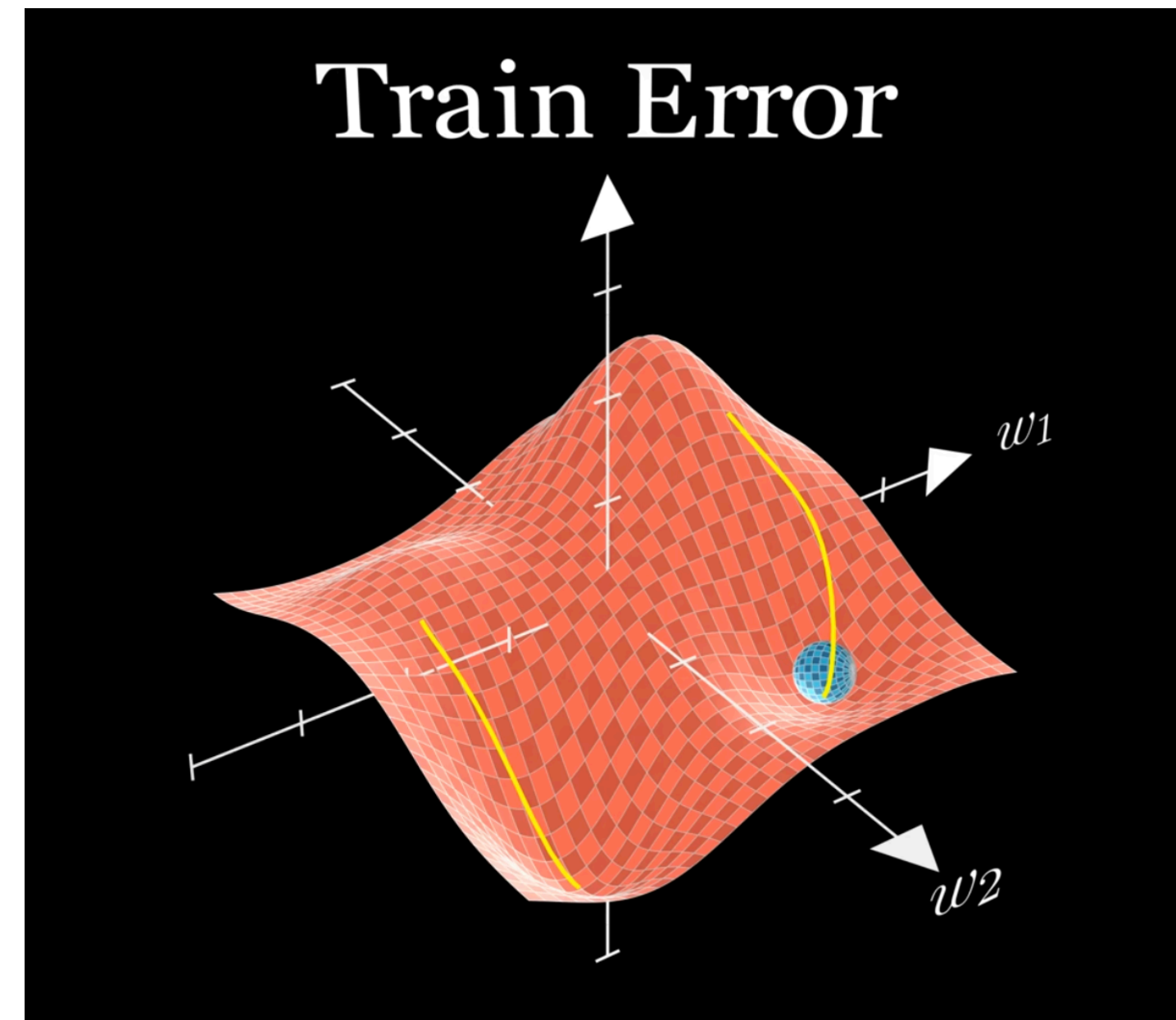
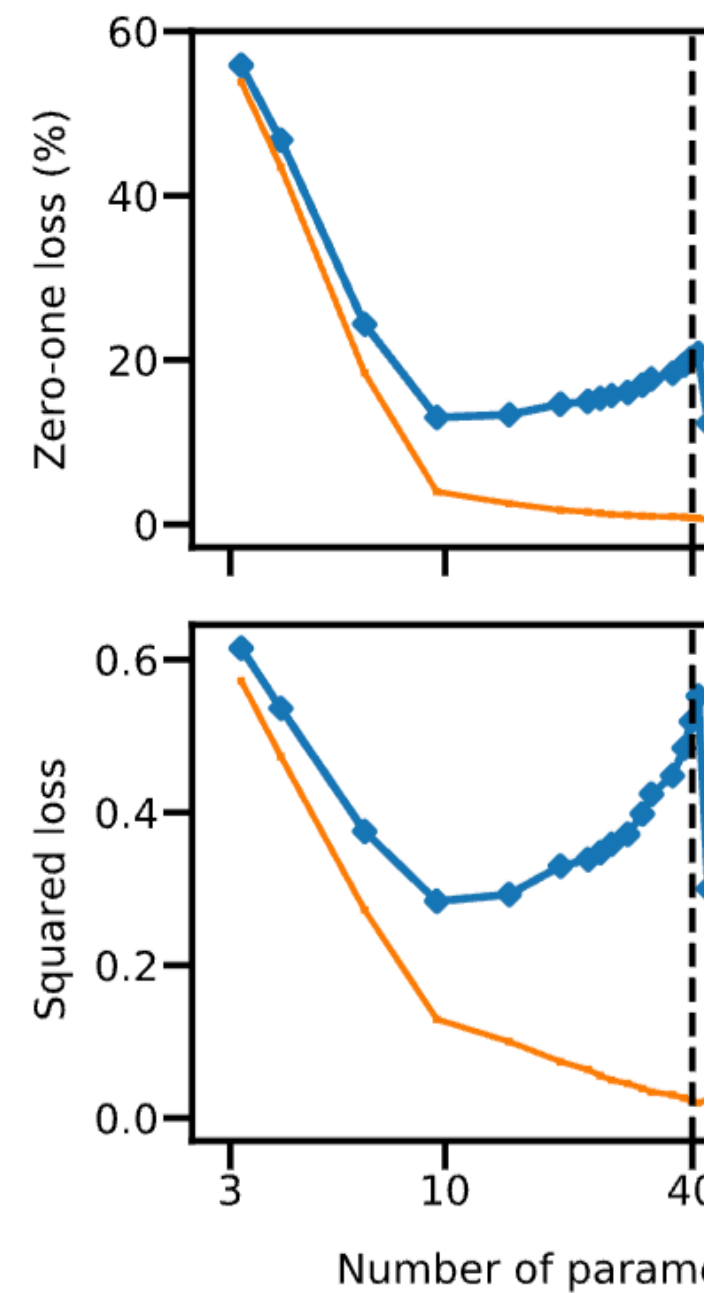
mcarbin@csail.mit.edu




The lottery ticket hypothesis:

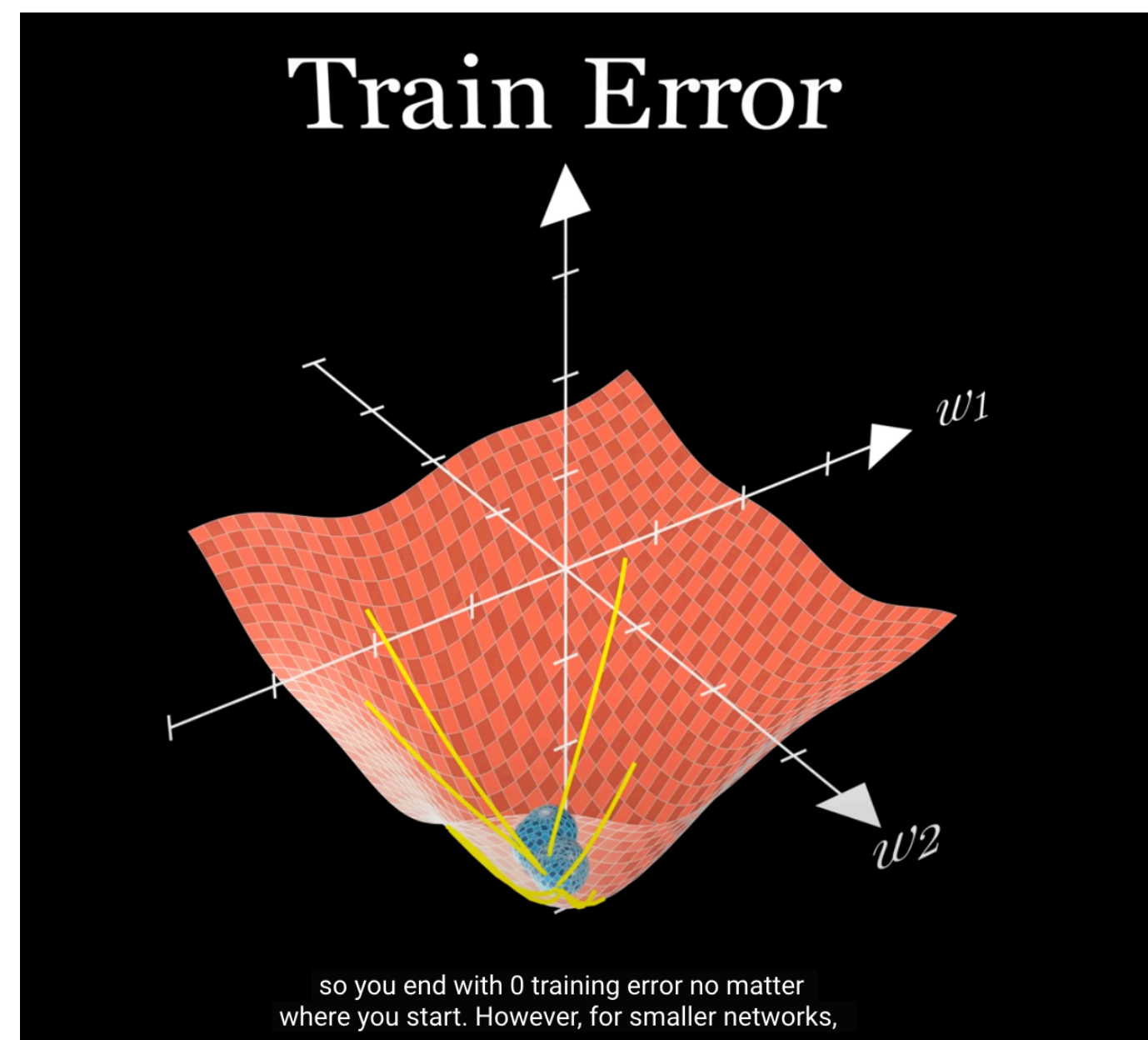
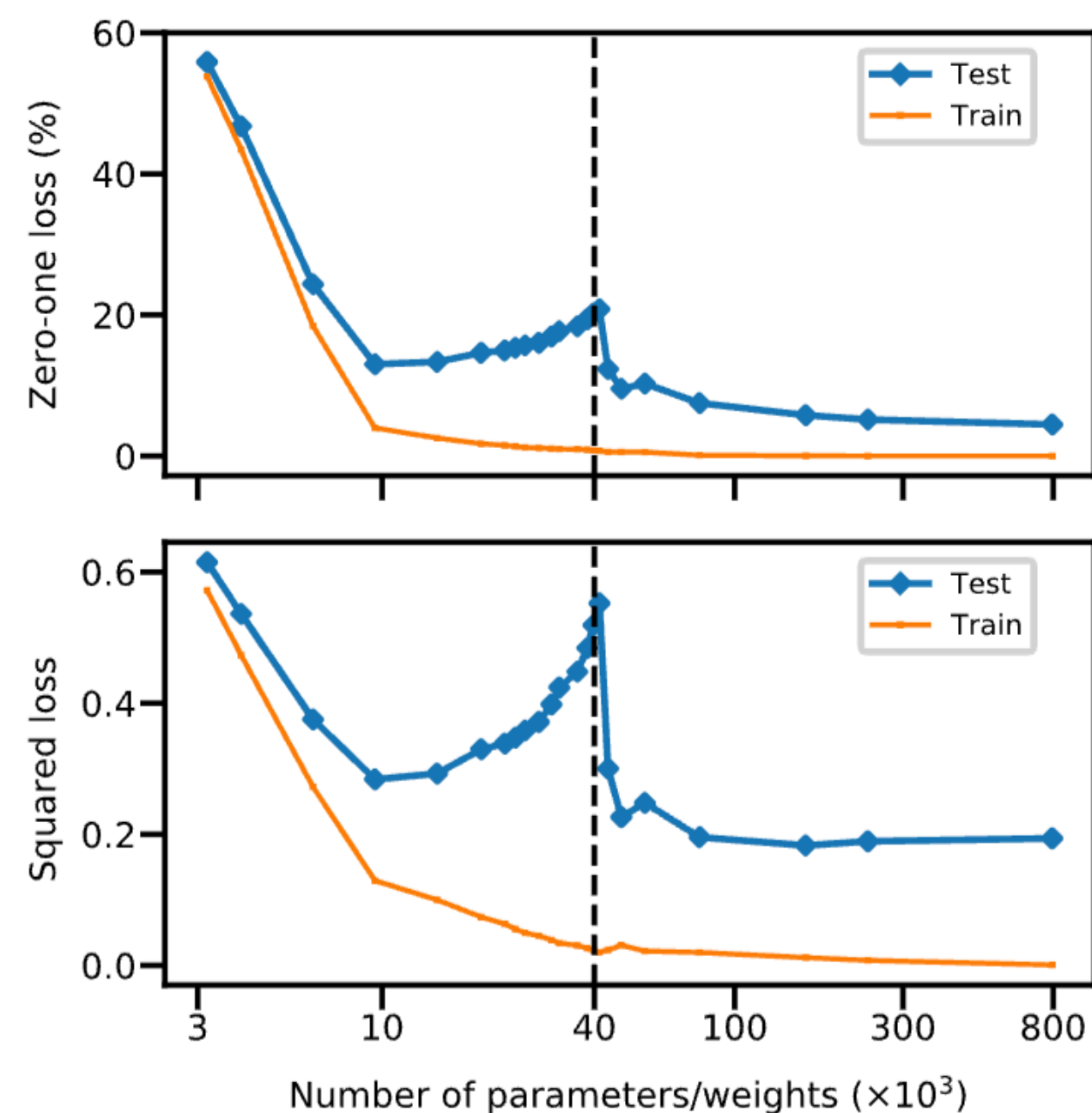
randomly-initialized NN contain **subnetworks** (winning tickets) can reach test accuracy **comparable to the original network** in a similar number of iterations.

The winning tickets we find have won the initialization lottery: their connections have initial weights that make training particularly effective.

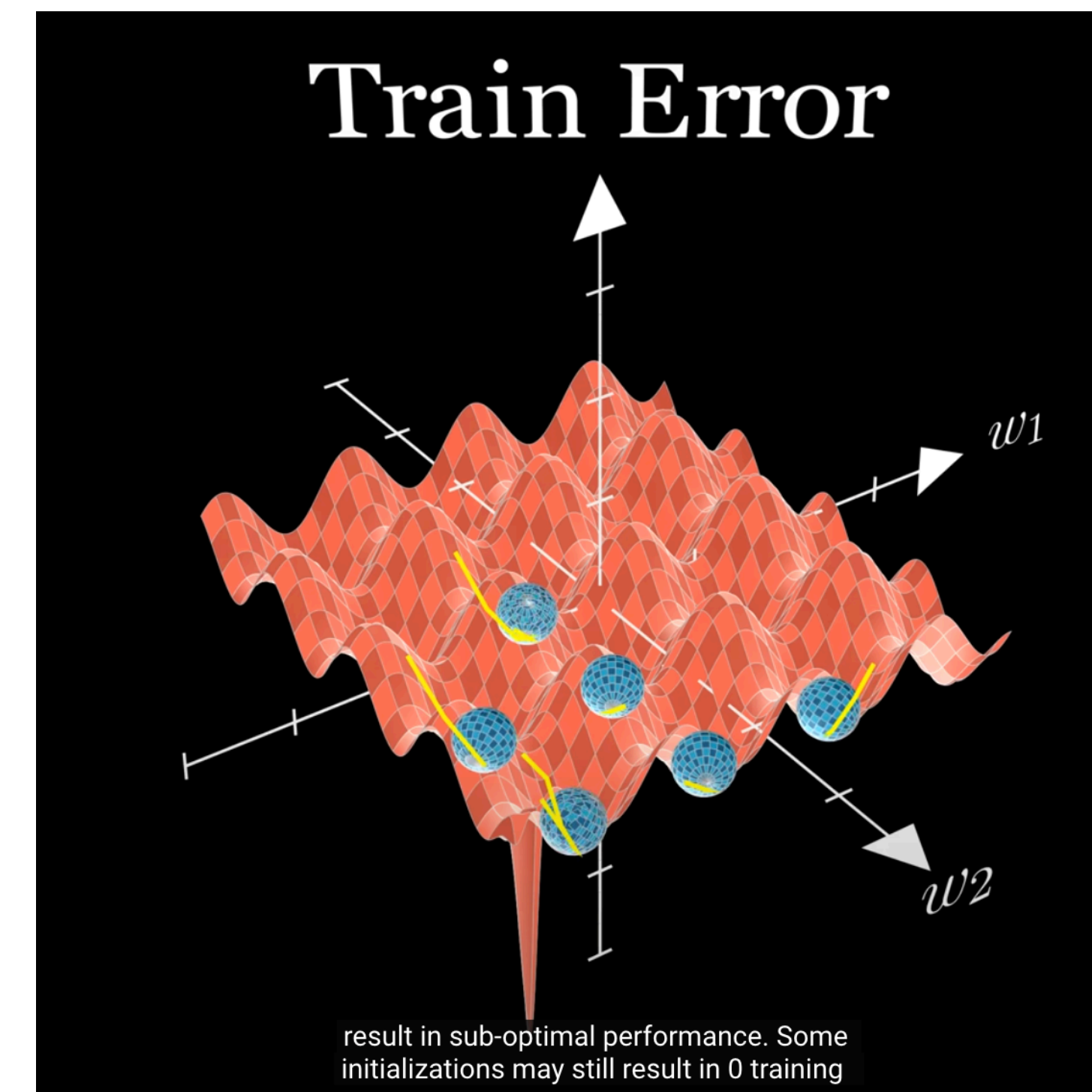


Receive ~ 0 training error no matter where you start.

However, for smaller networks, the initialization becomes crucial, as more possible initializations result in sub-optimal performances. (but s exist)



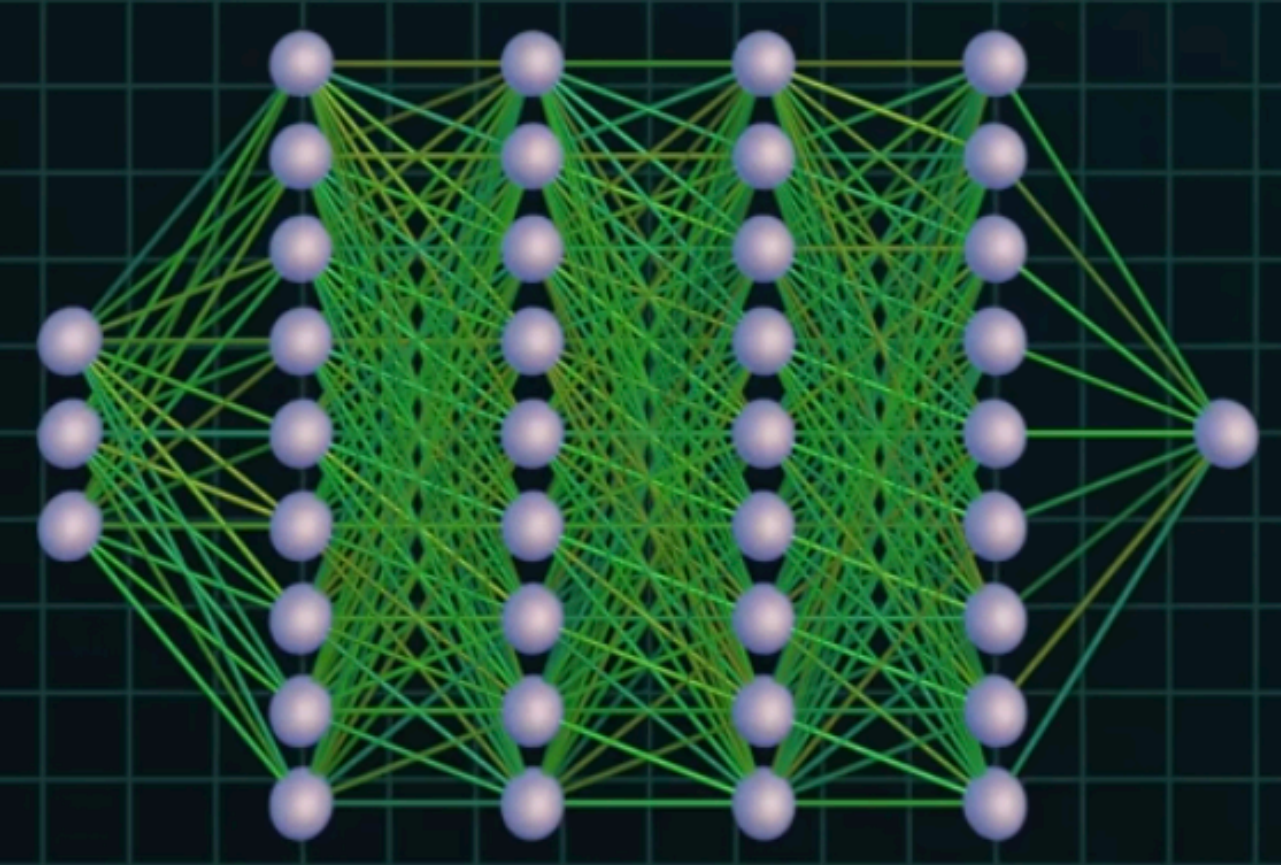
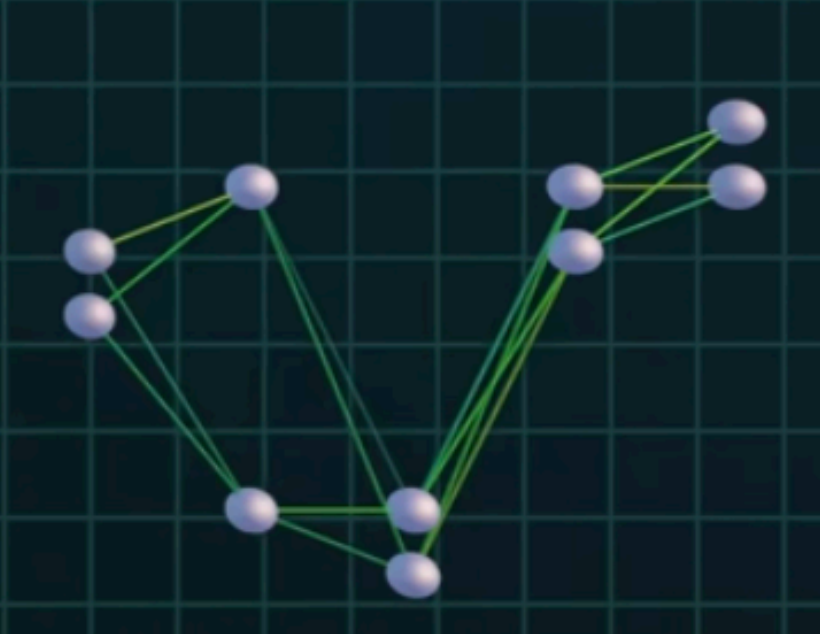
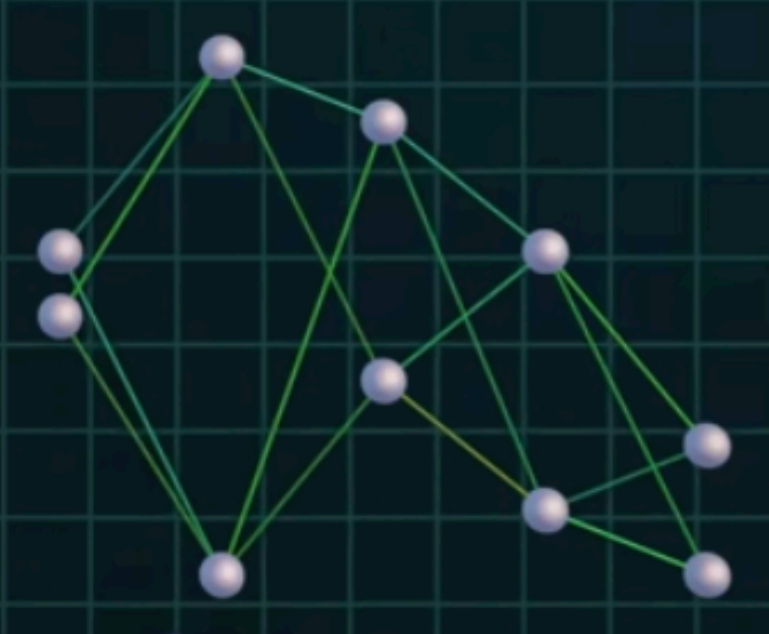
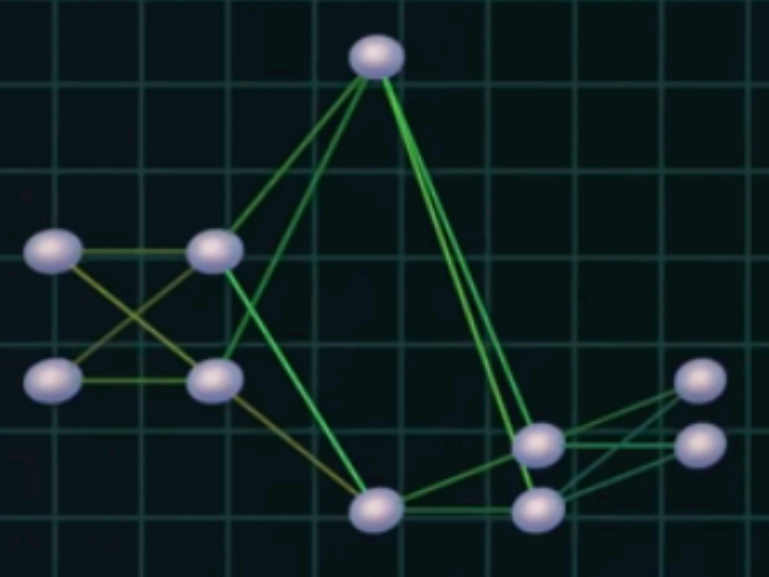
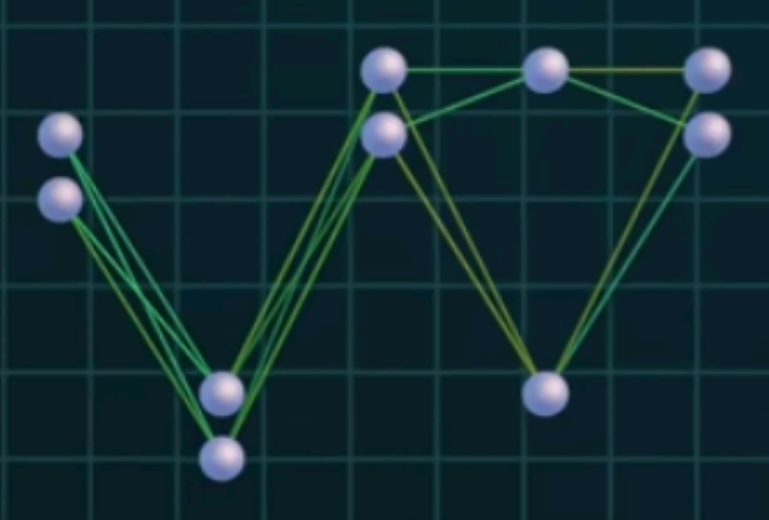
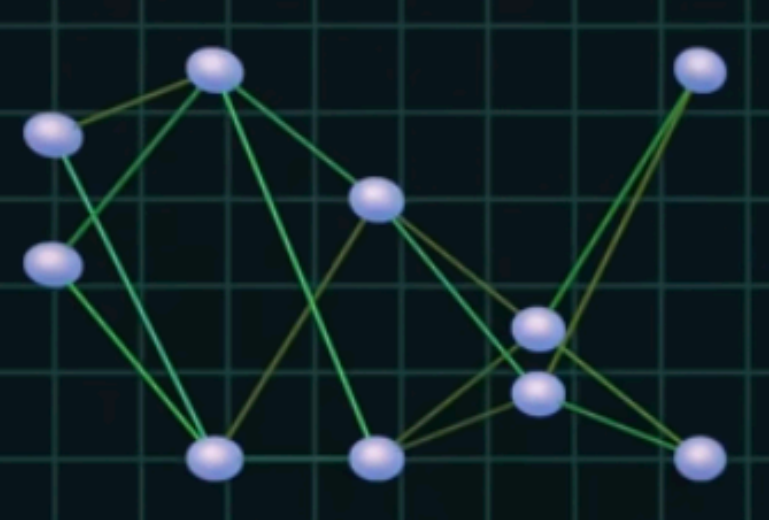
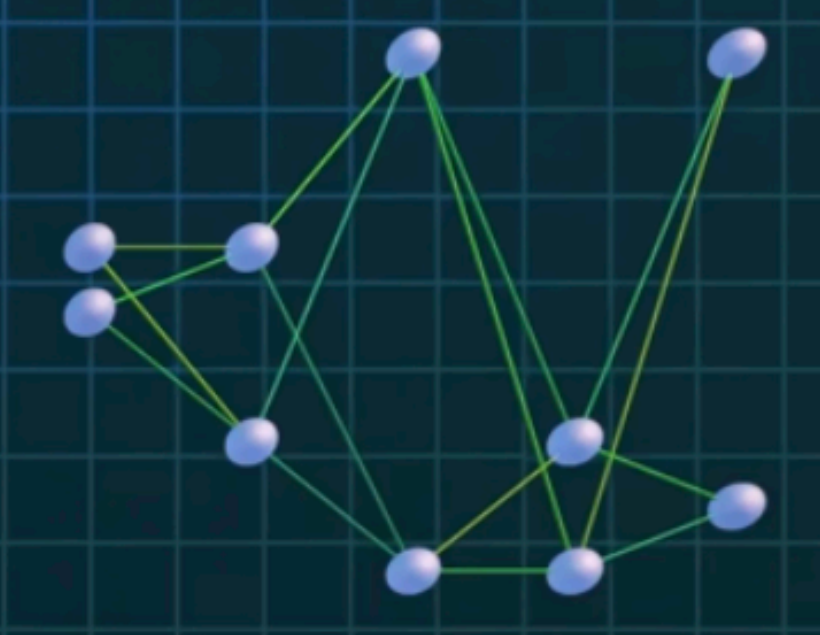
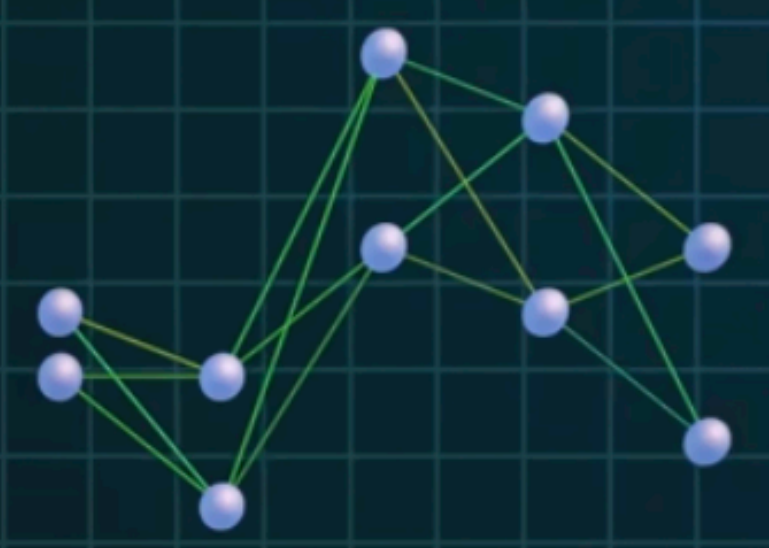
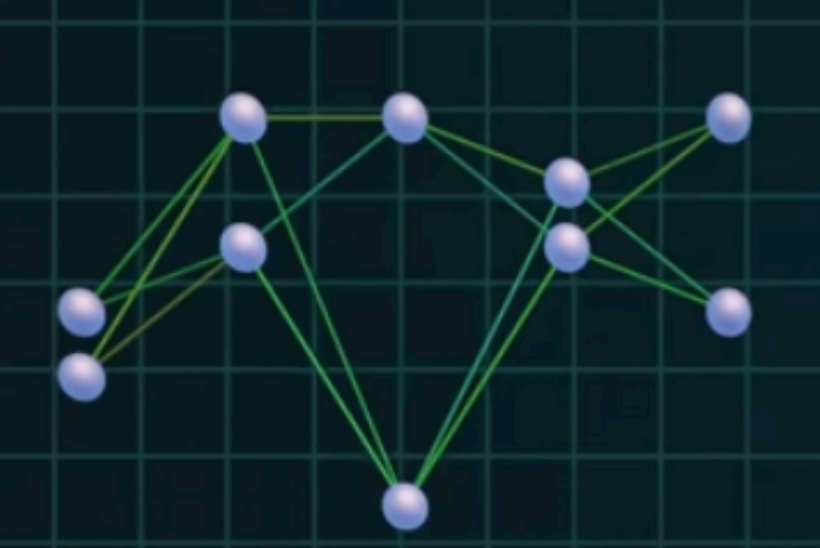
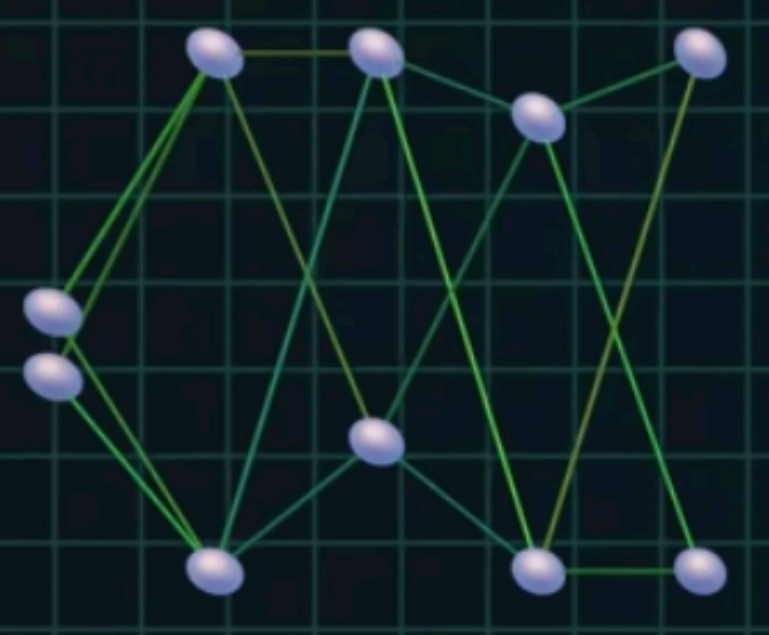
A large network



Many small sub-networks

Receive ~ 0 training error no matter where you start.

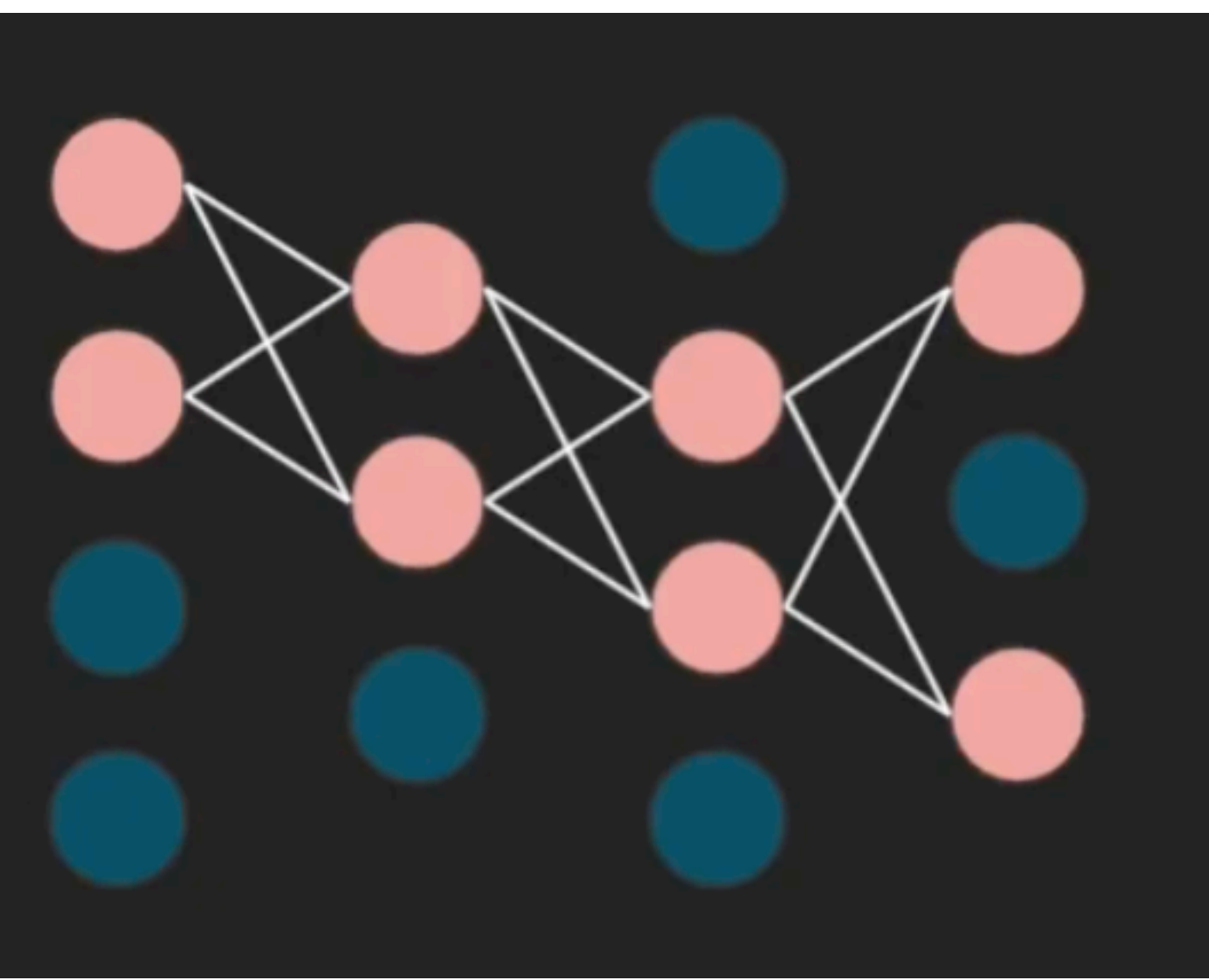
However, for smaller networks, the initialization becomes crucial, as **more possible initializations** result in sub-optimal performances. (but s **exist** in some rare initializations)



Hidden inside of a large neural net, is a much smaller **sub-network** that is actually doing all of the work. The rest is just useless fluff.

But why ?


why does it generalize better than a small network of the same size as this subnetwork.



The Winning Ticket(s)

Two ways to find a sub-network:

pruning procedure —> **uncovering the sub-network** with the better initialization.

larger network —> more s to choose from to win the lottery
—> **uncovering the sub-network** with the better initialization.

Pruning procedure =

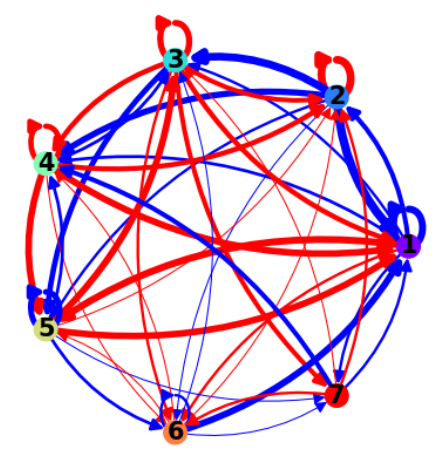
they were essentially uncovering the sub-network with the best initialization.

Questions to take away

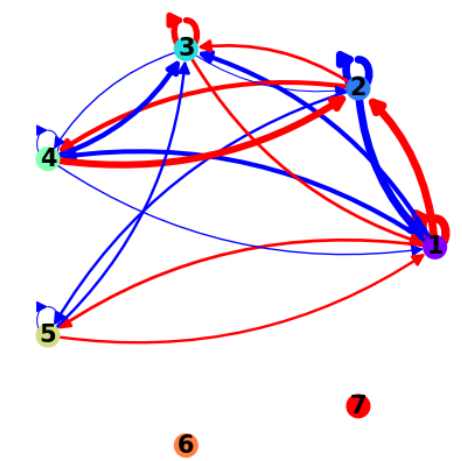
- fitness landscape
- GRNs are diverse
- GFN finds diversity

d Somite Patterns

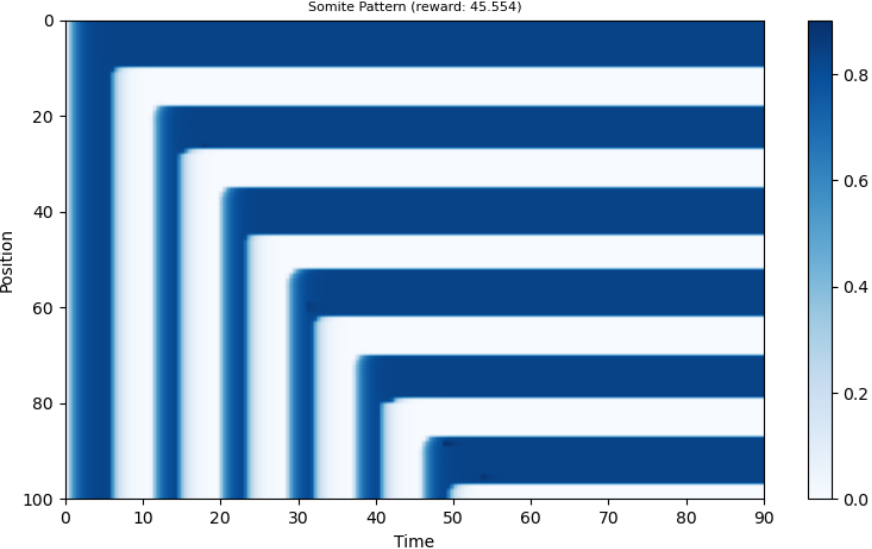
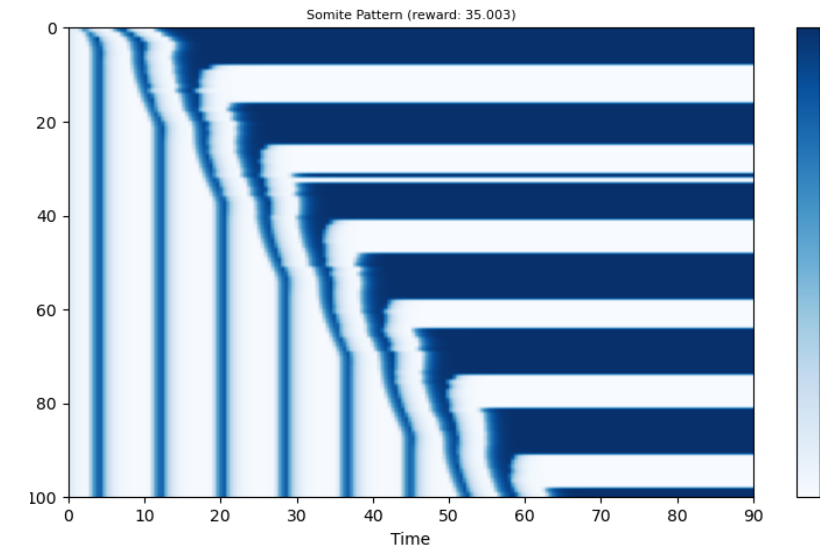
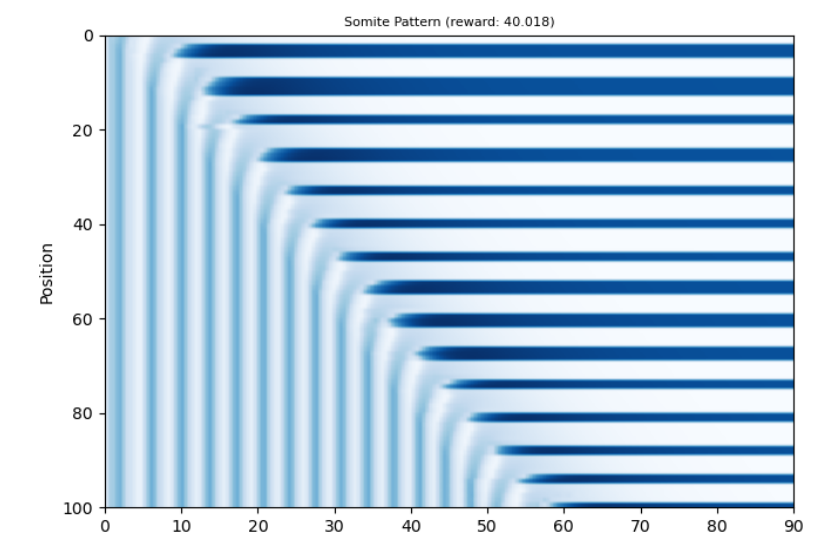
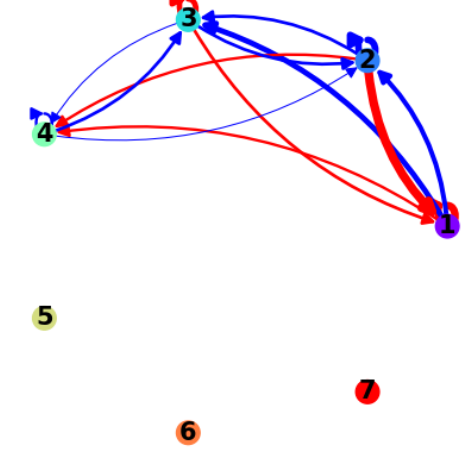
Motif 4									
-80	-50	-25	-30	80	25	-25	0	0	0
-100	90	-90	-70	-25	-5	0	0	0	0
55	90	70	55	60	25	50	0	0	0
75	55	-30	60	75	5	0	0	0	0
75	5	75	-30	65	-30	-5	0	0	0
-75	5	-5	5	5	-25	-5	0	0	0
-25	25	5	-50	0	25	0	0	0	0



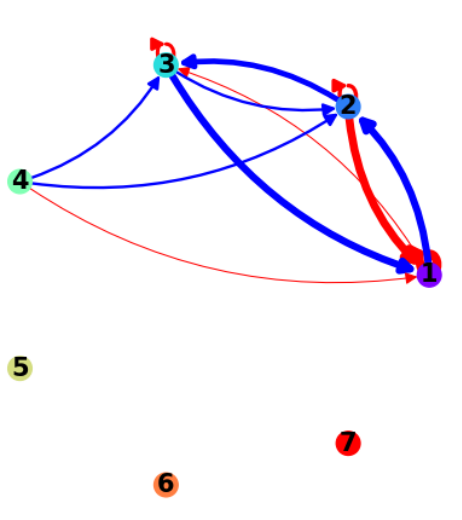
Motif 8									
100	90	-50	-50	30	0	0	0	0	0
-90	-90	30	50	-25	0	0	0	0	0
30	-5	70	-10	0	0	0	0	0	0
-10	80	-55	-5	0	0	0	0	0	0
25	0	-25	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0



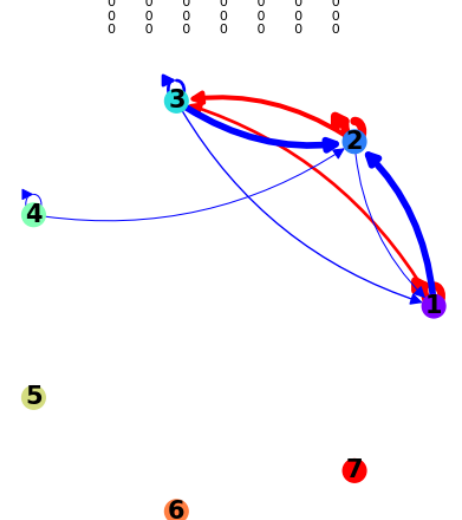
Motif 2									
85	-50	-60	25	0	0	0	0	0	0
100	-75	-30	0	0	0	0	0	0	0
30	-30	60	-5	0	0	0	0	0	0
0	-5	-30	-25	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0



Motif 22									
100	-80	5	0	0	0	0	0	0	0
90	40	-60	0	0	0	0	0	0	0
-80	-25	30	0	0	0	0	0	0	0
5	-25	-25	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0



Motif 1									
90	-75	30	0	0	0	0	0	0	0
100	50	0	0	0	0	0	0	0	0
-75	-30	0	0	0	0	0	0	0	0
-5	-5	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0



Motif 2									
85	-80	0	0	0	0	0	0	0	0
10	50	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0

