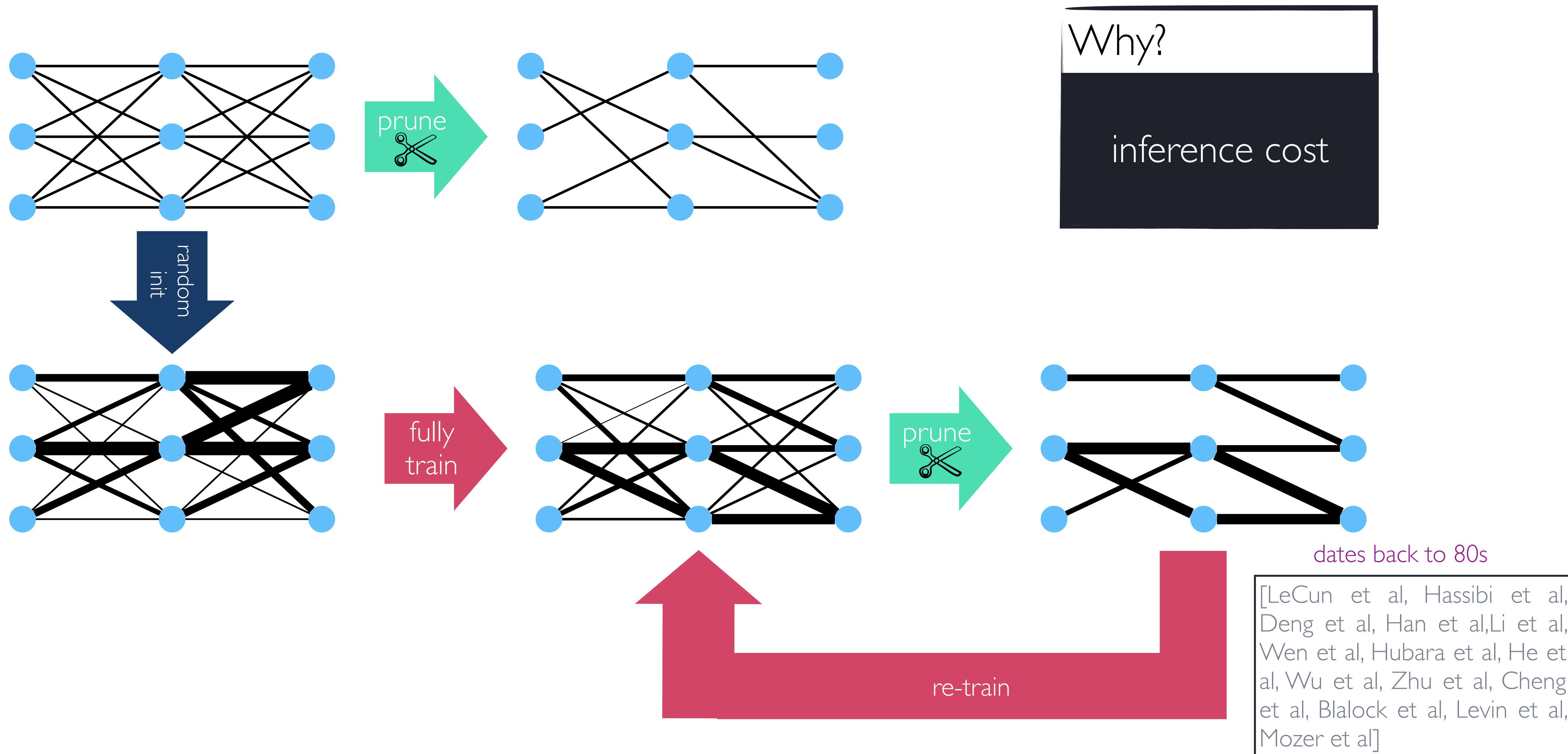# From Model Pruning to Sparse Updates and the Lottery Ticket Hypothesis
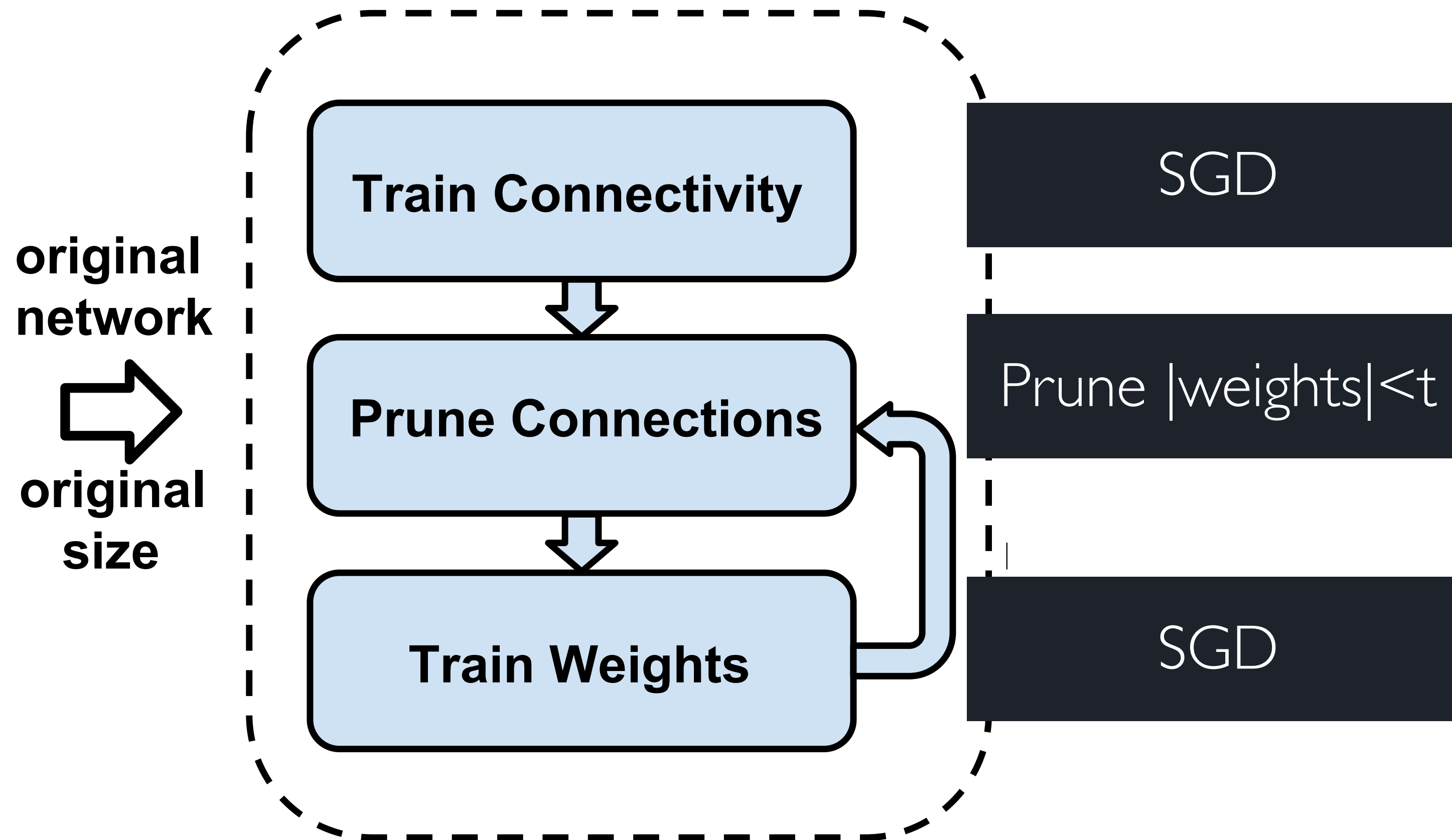
# Network Pruning, 1980-2018
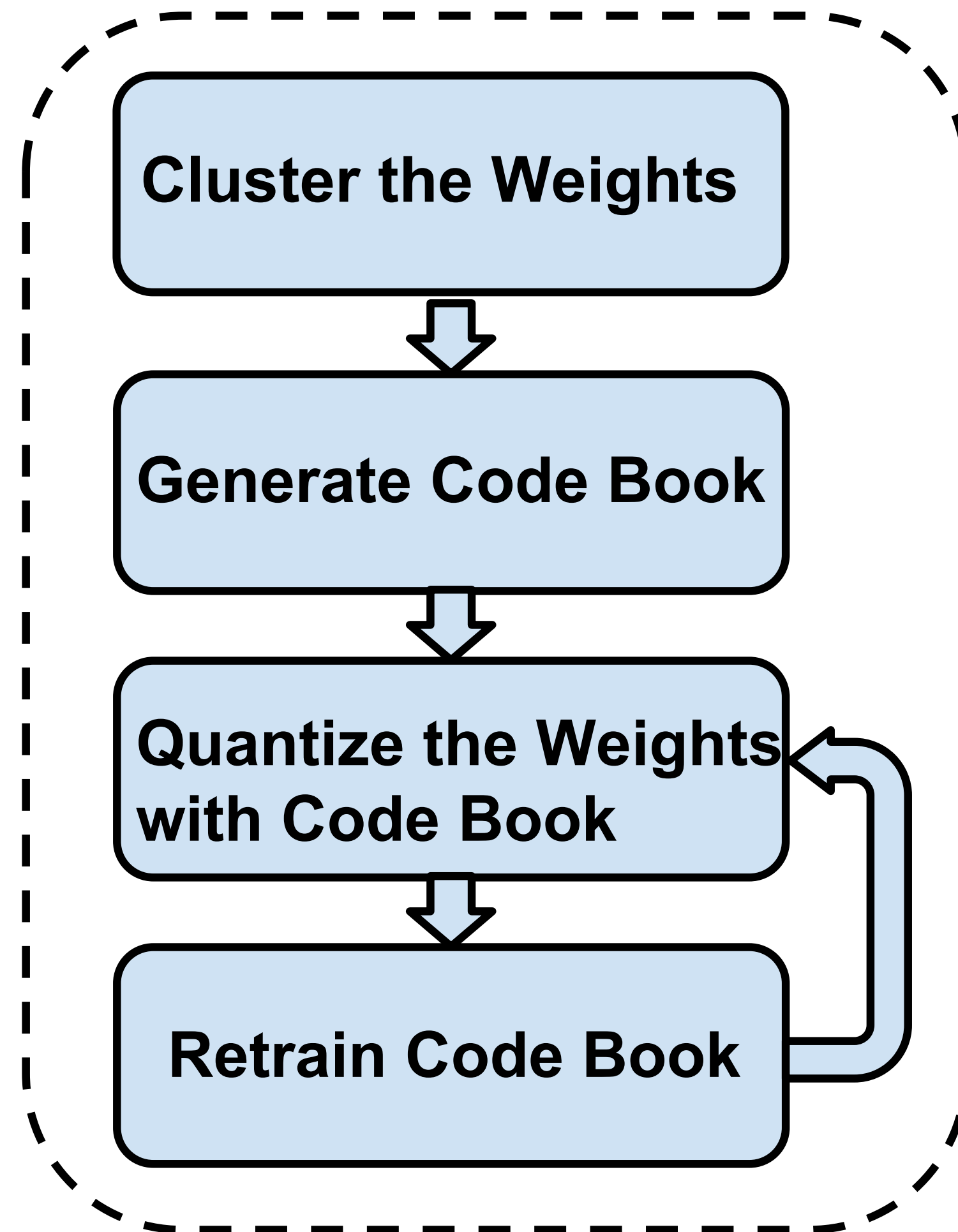


prune

random init

fully train

prune

re-train

Why?

inference cost

dates back to 80s

[LeCun et al, Hassibi et al, Deng et al, Han et al, Li et al, Wen et al, Hubara et al, He et al, Wu et al, Zhu et al, Cheng et al, Blalock et al, Levin et al, Mozer et al]

# An example: Deep Compression [ICLR, 2016]

**Quantization: less bits per weight**

**Pruning: less number of weights**

**Huffman Encoding**

original network

original size

**Train Connectivity**

**Prune Connections**

**Train Weights**

same accuracy

9x-13x reduction

**Cluster the Weights**

**Generate Code Book**

**Quantize the Weights with Code Book**

**Retrain Code Book**

same accuracy

27x-31x reduction

**Encode Weights**

**Encode Index**

same accuracy

35x-49x reduction

[Han, Mao, Dally, ICLR 2016]

# Deep Compression: Step 1, prune

**Train Connectivity**

SGD

original network ⇒ original size

**Prune Connections**

Prune |weights|<t

**Train Weights**

SGD

[Han, Mao, Dally, ICLR 2016]

# Deep Compression: Step 2, quantize



$$\min_{S_1,...,S_k} \sum_{i=1}^{k} \sum_{x \in S_i} \|x - \mu_i\|^2$$

K-means used
for clustering

[Han, Mao, Dally, ICLR 2016]

# Deep Compression: Step 3, compress

**Huffman Encoding**

# Deep Compression: Experiments

| Network | Top-1 Error | Top-5 Error | Parameters | Compress Rate |
|---|---|---|---|---|
| LeNet-300-100 Ref | 1.64% | - | 1070 KB | |
| LeNet-300-100 Compressed | 1.58% | - | **27 KB** | 40× |
| LeNet-5 Ref | 0.80% | - | 1720 KB | |
| LeNet-5 Compressed | 0.74% | - | **44 KB** | **39×** |
| AlexNet Ref | 42.78% | 19.73% | 240 MB | |
| AlexNet Compressed | 42.78% | 19.70% | **6.9 MB** | **35×** |
| VGG-16 Ref | 31.50% | 11.32% | 552 MB | |
| VGG-16 Compressed | 31.17% | 10.91% | **11.3 MB** | 49× |

network pruning works!

# Network Pruning, 1980-2018

prune

Why?

inference cost

random init

## Issue:
training to full acc and THEN pruning
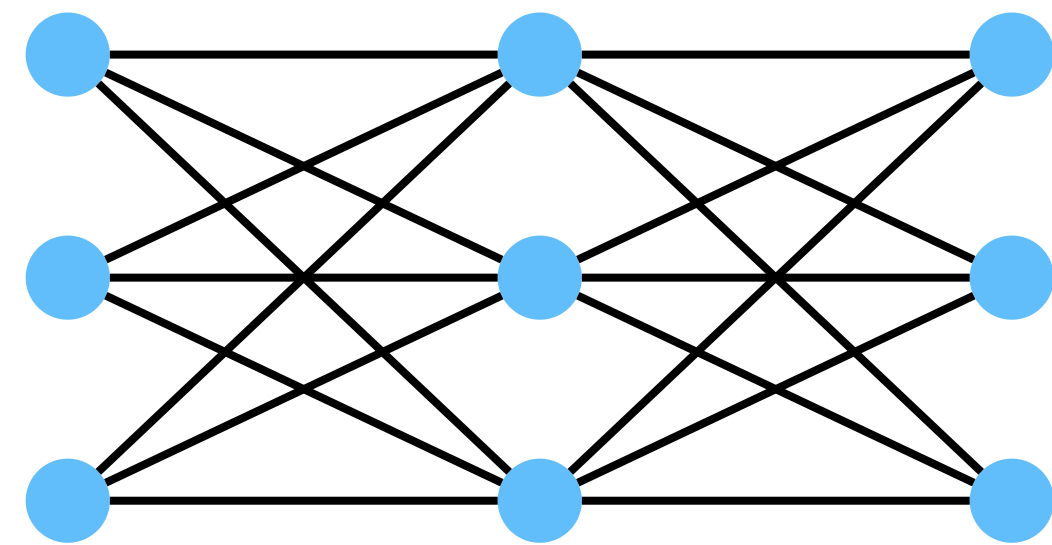is *expensive*
Q: Can we avoid?

re-train

dates back to 80s

[LeCun et al, Hassibi et al, Deng et al, Han et al, Li et al, Wen et al, Hubara et al, He et al, Wu et al, Zhu et al, Cheng et al, Blalock et al, Levin et al, Mozer et al]
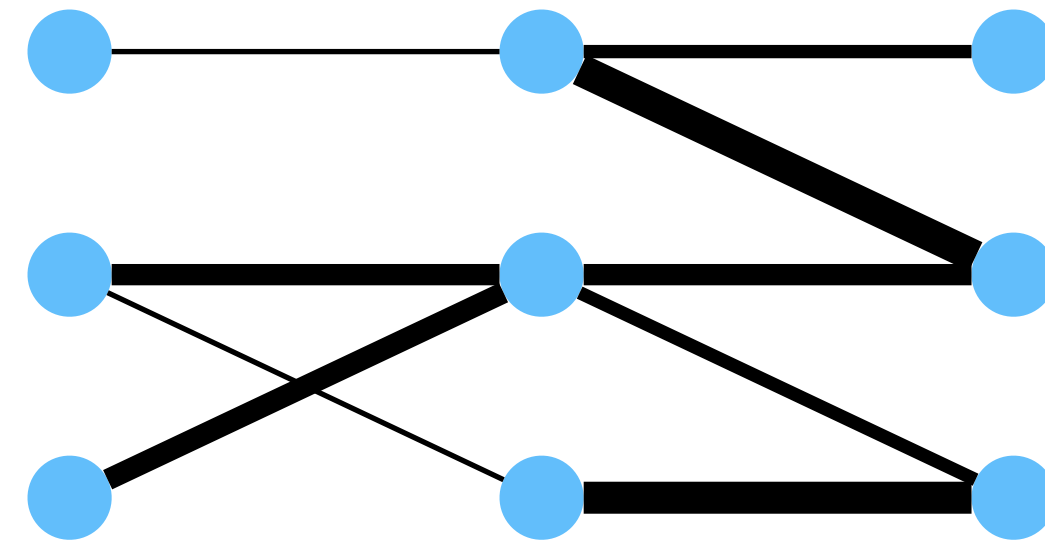
# The Lottery Ticket Hypothesis
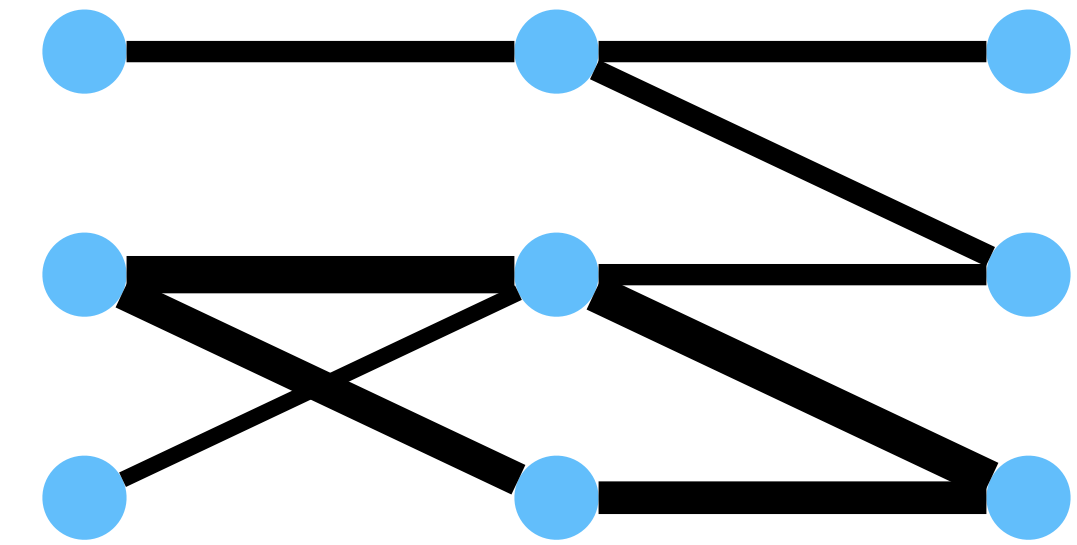
# Lottery Ticket Hypothesis (LTH)
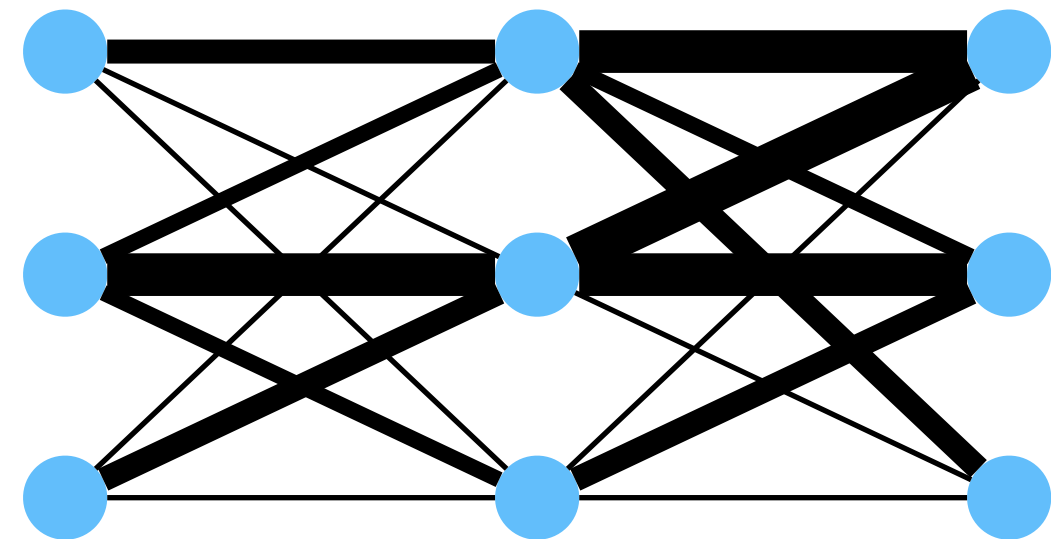
Frankle, Carbin, ICLR 2019



win the lottery

train

random init

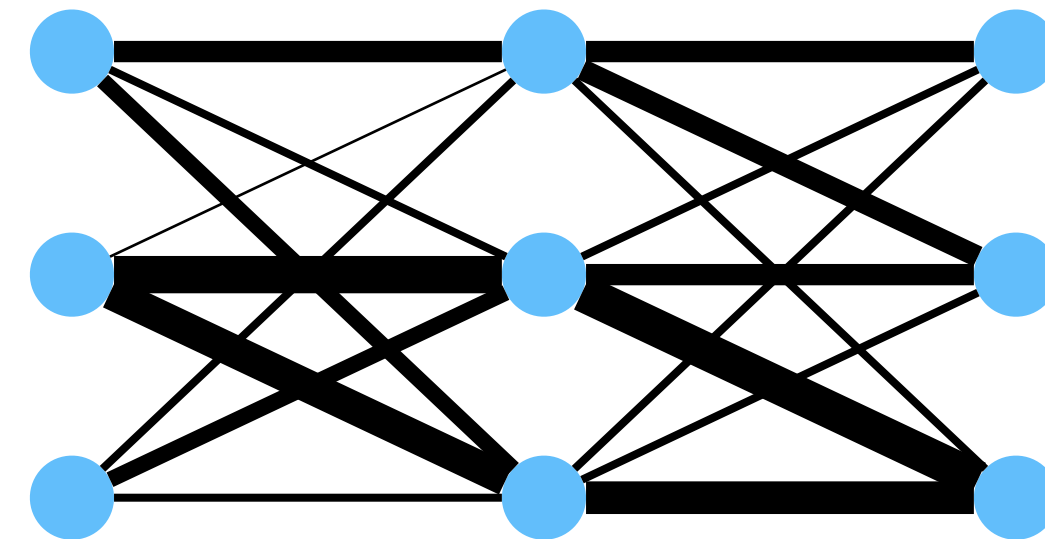"there exist sparse subnets at init trainable to full accuracy"
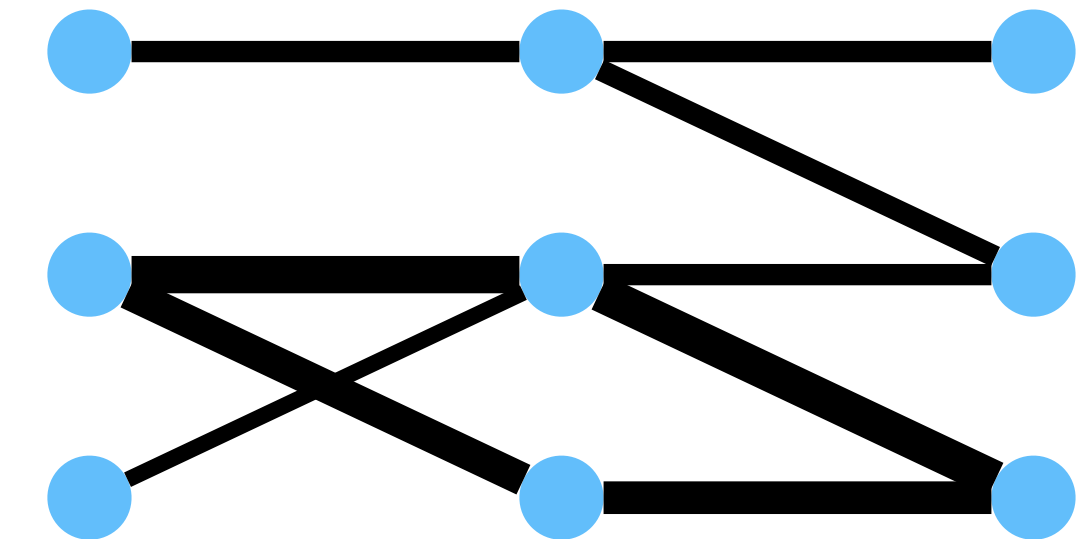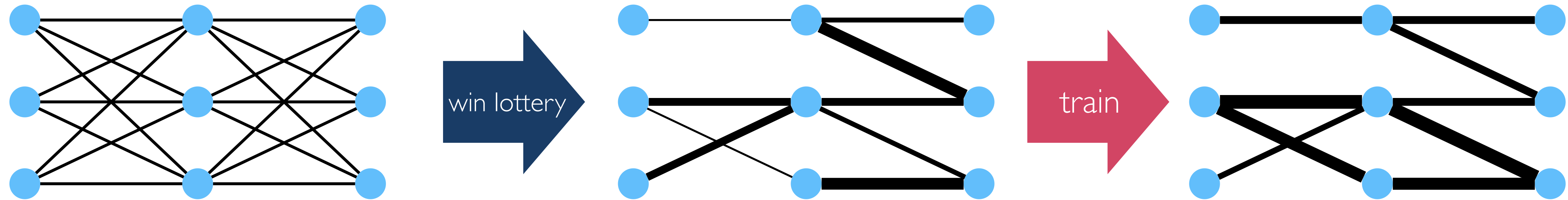
train

prune

re-train

# Lottery Ticket Hypothesis (LTH)

Frankle, Carbin, ICLR 2019



"there exist sparse subnetworks at init that can be trained to full accuracy"
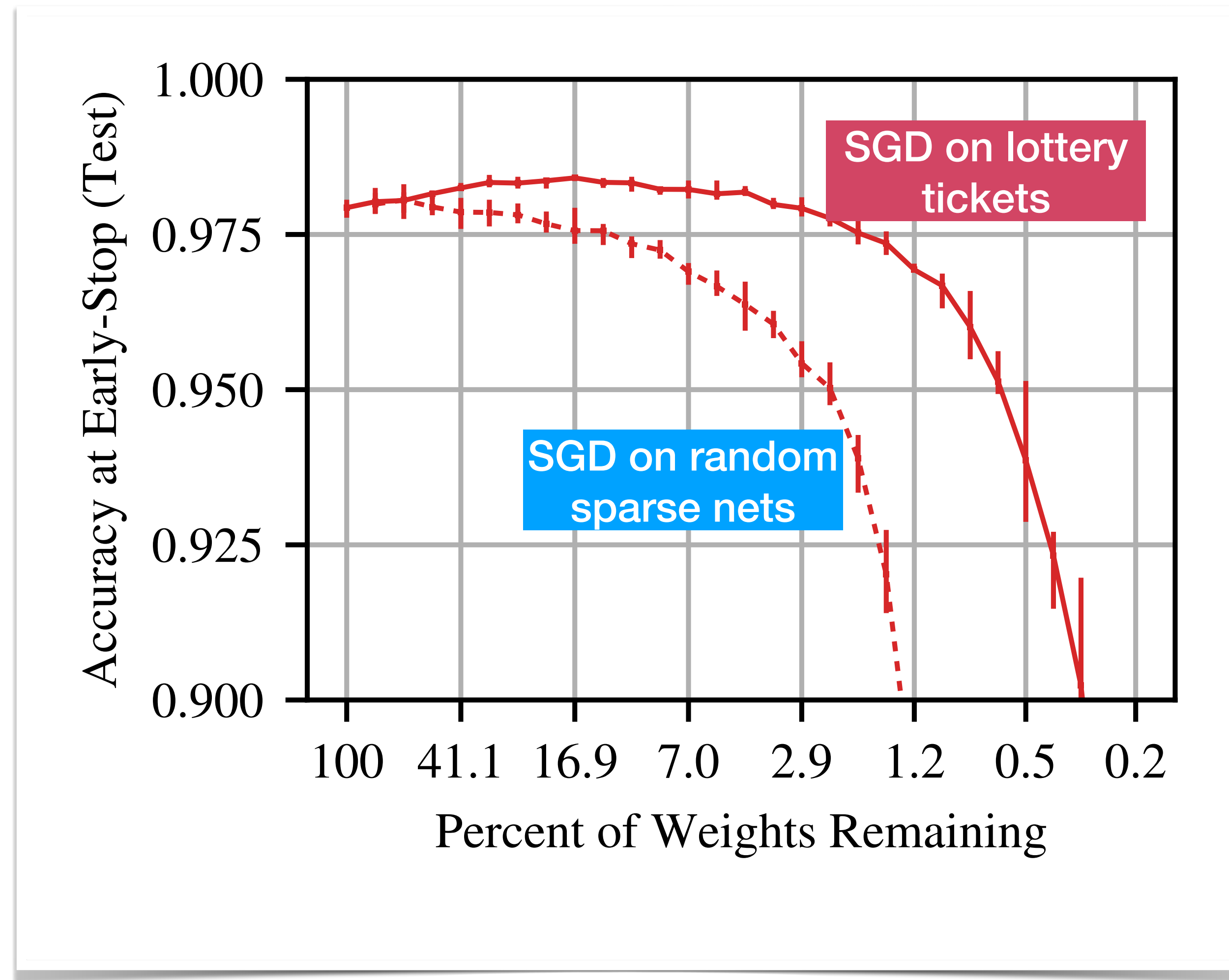
Identify

If true, kinda big deal!
We can avoid pruning/retraining cycle
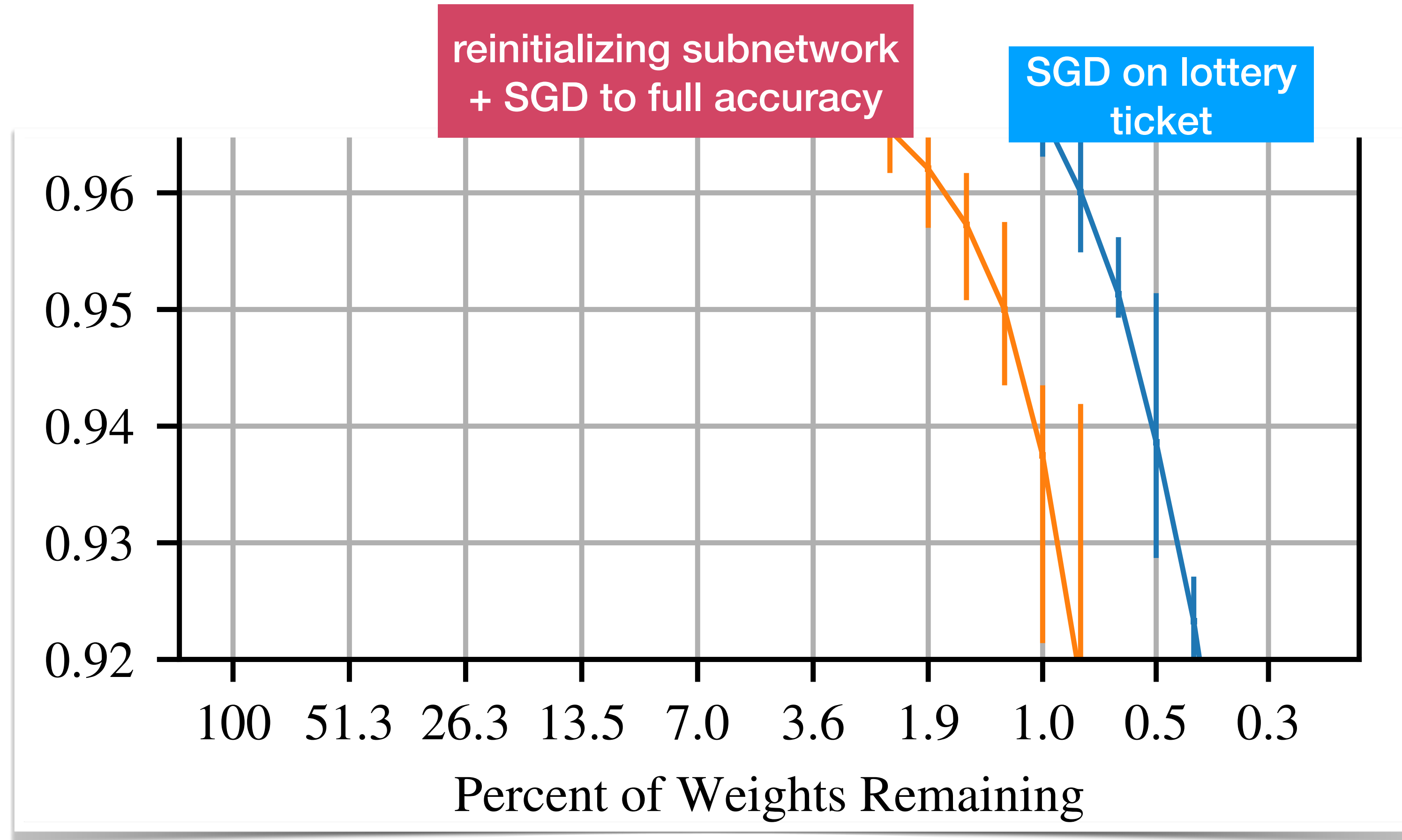

Q: How do you win the lottery??

5. (Sometimes) GOTO 3

note: winning the lottery does
not come for free

# Lottery tickets >> random subnets

# lottery ticket = subnetwork + original weights



reinitializing subnetwork + SGD to full accuracy

SGD on lottery ticket

0.96
0.95
0.94
0.93
0.92

100  51.3  26.3  13.5  7.0  3.6  1.9  1.0  0.5  0.3

Percent of Weights Remaining

LTH research has been very active:
[Zhou, Lan, Liu, Yosinski]
[Frankle, Dziugaite, Roy, Carbin]
[Cosentino, Zaiter, Pei, Zhu]
[Soelen, Sheppard]
[Sabatelli, Kestemont, Geurts]
[Ramanujan, Wortsman, Kembhavi, Farhadi, Rastegari]
[Wang, Zhang, Xie, Zhou, Su, Zhang, Hu]

# Many many extensions

## The Lottery Ticket Hypothesis for Pre-trained BERT Networks

Tianlong Chen[1], Jonathan Frankle[2], Shiyu Chang[3], Sijia Liu[3], Yang Zhang[3],
Zhangyang Wang[1], Michael Carbin[2]
[1]University of Texas at Austin, [2]MIT CSAIL, [3]MIT-IBM Watson AI Lab, IBM Research
{tianlong.chen,atlaswang}@utexas.edu,{jfrankle,mcarbin}@csail.mit.edu,
{shiyu.chang,sijia.liu,yang.zhang2}@ibm.com

## One ticket to win them all: generalizing lottery ticket initializations across datasets and optimizers

**Ari S. Morcos***
Facebook AI Research
arimorcos@fb.com

**Haonan Yu**
Facebook AI Research
haonanu@gmail.com

**Michela Paganini**
Facebook AI Research
michela@fb.com

**Yuandong Tian**
Facebook AI Research
yuandong@fb.com

## One ticket to win them all: generalizing lottery ticket initializations across datasets and optimizers

**Ari S. Morcos***
Facebook AI Research
arimorcos@fb.com

**Haonan Yu**
Facebook AI Research
haonanu@gmail.com

**Michela Paganini**
Facebook AI Research
michela@fb.com

**Yuandong Tian**
Facebook AI Research
yuandong@fb.com

## DRAWING EARLY-BIRD TICKETS: TOWARDS MORE EFFICIENT TRAINING OF DEEP NETWORKS

**Haoran You, Chaojian Li, Pengfei Xu, Yonggan Fu, Yue Wang, Richard G. Baraniuk & Yingyan Lin***
Department of Electrical and Computer Engineering
Rice University
Houston, TX 77005, USA
{hy34, cl114, px5, yf22, yw68, yingyan.lin, richb}@rice.edu

**Xiaohan Chen & Zhangyang Wang***
Department of Computer Science and Engineering
Texas A&M University
College Station, TX 77843, USA
{chernxh, atlaswang}@tamu.edu

## Rigging the Lottery: Making All Tickets Winners

**Utku Evci**[1]  **Trevor Gale**[1]  **Jacob Menick**[2]  **Pablo Samuel Castro**[1]  **Erich Elsen**[2]

## PUFFERFISH: COMMUNICATION-EFFICIENT MODELS AT NO EXTRA COST

**Hongyi Wang,**[1]  **Saurabh Agarwal,**[1]  **Dimitris Papailiopoulos**[2]

# Challenges of Pruning At Initialization

# Pruning at initialization does't seem work

RETHINKING VALUE OF NETWORKS PRUNING

Zhuang Liu[1]*, Mingjie Sun[2]*†, Tinghui Zhou[1], Gao Huang[2], Trevor Darrell
[1]University of California, Berkeley   [2]Tsinghua University

Random networks areas good as (or better than) claimed lottery tickets..

PRUNING NEURAL NETWORKS AT INITIALIZATION: WHY ARE WE MISSING THE MARK?

Jonathan Frankle    Gintare Karolina Dziugaite    Daniel M. Roy    Michael Carbin
MIT CSAIL           Element AI                     University of Toronto   MIT CSAIL
                                                   Vector Institute
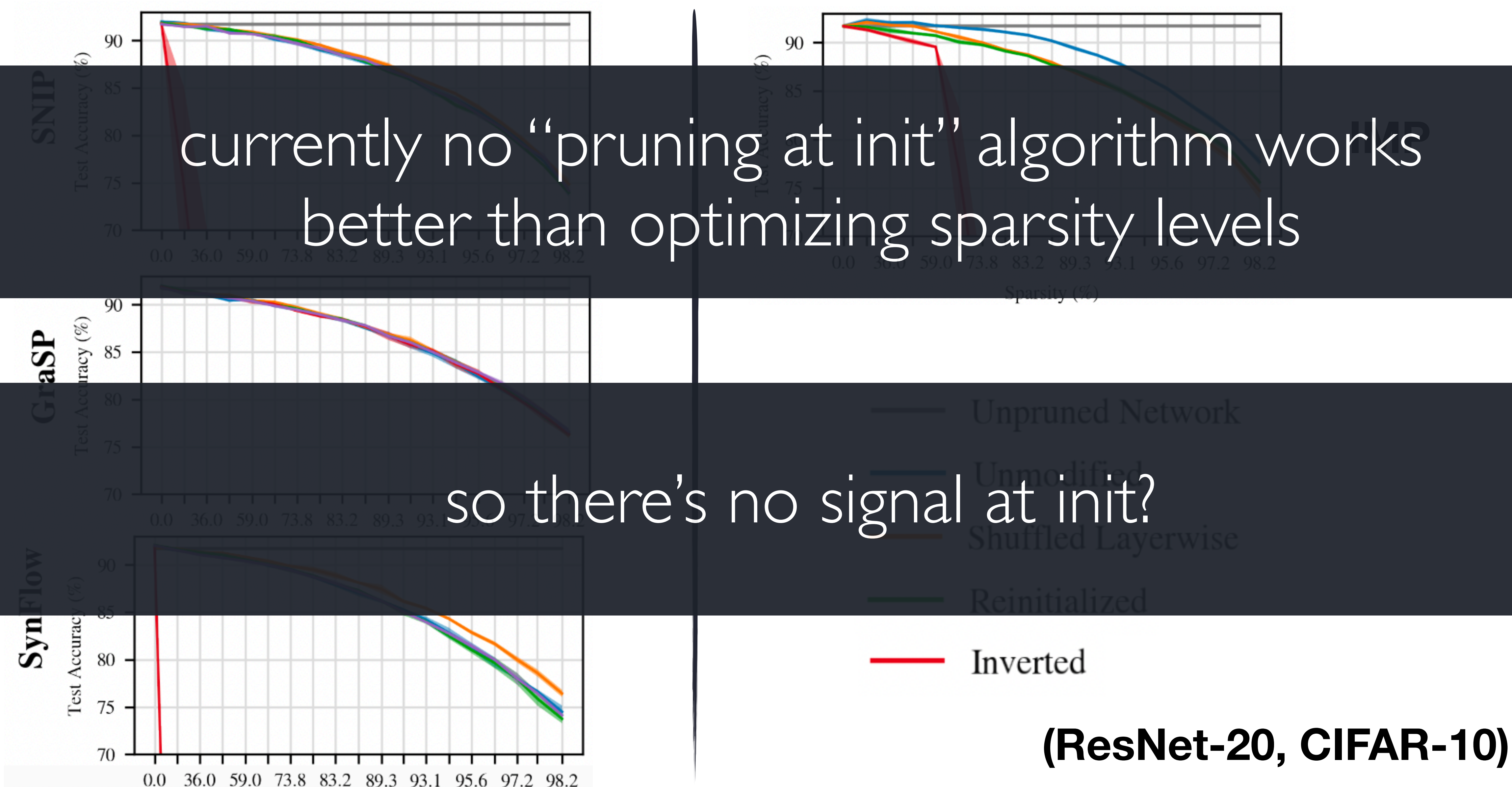
**Sanity-Checking Pruning Methods:**
**Random Tickets can Win the Jackpot**

THE EARLY PHASE OF NEURAL NETWORK TRAINING

Finding good tickets at init seems hard.
Current fix = wait for a few epochs

Jingtong Su[1]*   Yihang Chen[2]*   Tianle Cai[4]*   Jonathan Frankle   David Schwab   Ari S. Morcos
                                                    MIT CSAIL          CCN CUNY       Facebook AI Research
                                                                       Facebook AI Research

Tianhao Wu[2]   Ruiqi Gao[3,4]   Liwei Wang[5]   Jason D. Lee[3]

# Sanity Checks [Frankle'21]



currently no "pruning at init" algorithm works
better than optimizing sparsity levels

so there's no signal at init?

Unpruned Network

Unmodified

Shuffled Layerwise

Reinitialized

Inverted

**(ResNet-20, CIFAR-10)**
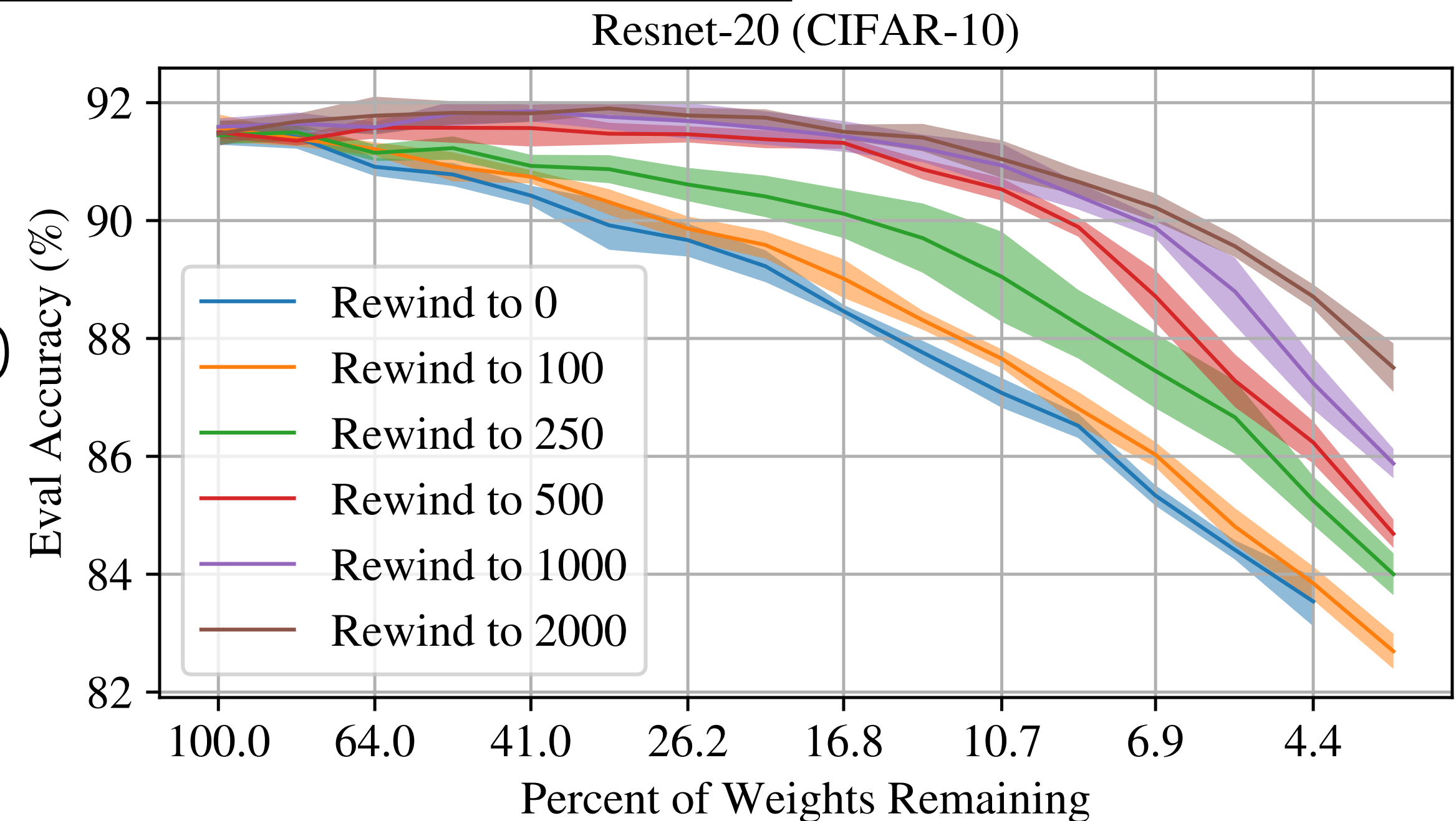
# Fixing IMP

# fix = rewind shortly after init

THE EARLY PHASE OF NEURAL NETWORK TRAINING

**Jonathan Frankle**[†]
MIT CSAIL

**David J. Schwab**
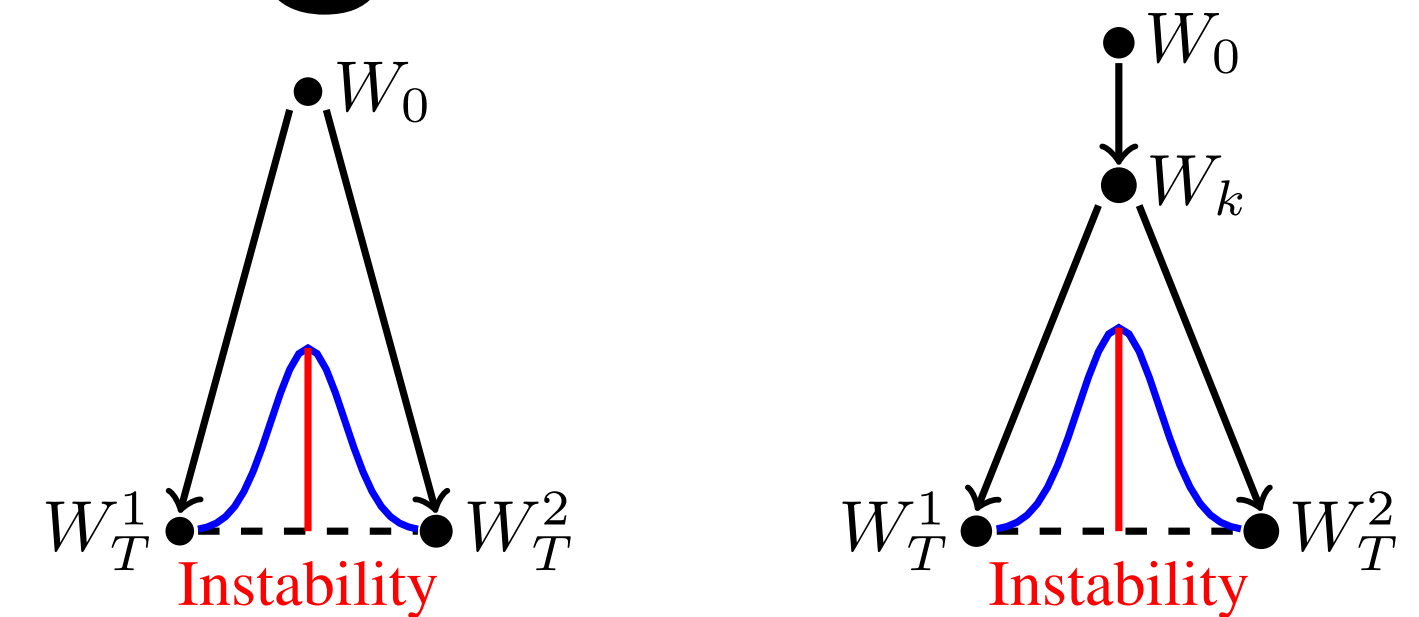CUNY ITS
Facebook AI Research

**Ari S. Morcos**
Facebook AI Research

- "vanilla" LTH only true for MNIST/small nets
- Rewinding to init doesn't work for Resnet/Cifar10
- One needs to rewind later (i.e., train a bit)

Resnet-20 (CIFAR-10)

# Lottery Tickets are hard to get at Init

**Linear Mode Connectivity and the Lottery Ticket Hypothesis**

Jonathan Frankle [1]   Gintare Karolina Dziugaite [2]   Daniel M. Roy [3,4]   Michael Carbin [1]
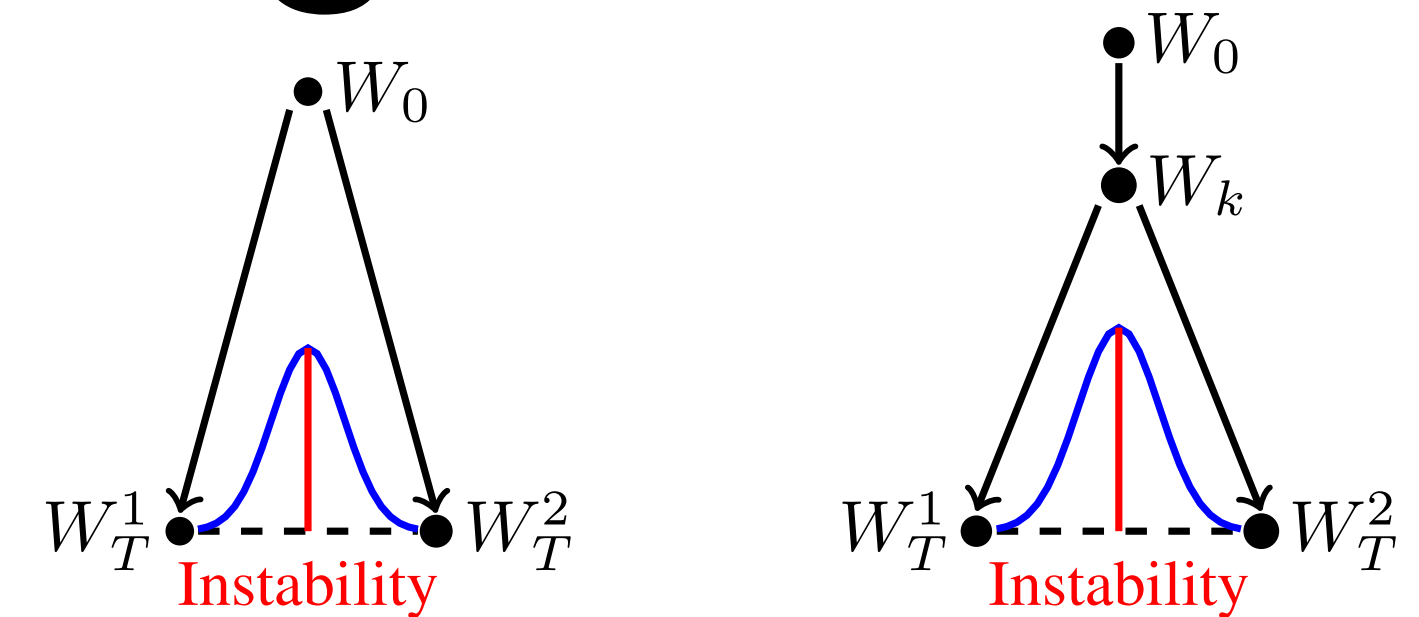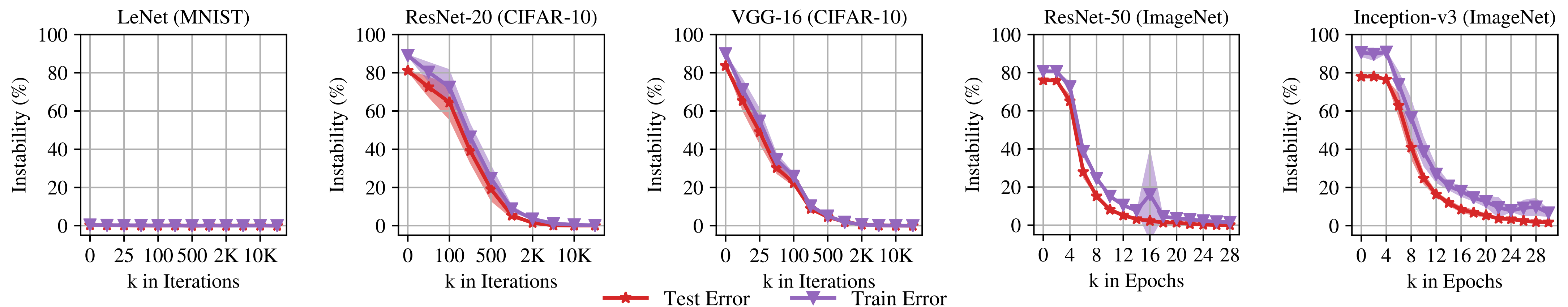
*Figure 1.* A diagram of instability analysis from step 0 (left) and step $k$ (right) when comparing networks using linear interpolation.

- Rewinding to iteration K, rather than init works much better
- Experimental analysis through the existence of linear connectivity

# Lottery Tickets are hard to get at Init



**Linear Mode Connectivity and the Lottery Ticket Hypothesis**

Jonathan Frankle[1]   Gintare Karolina Dziugaite[2]   Daniel M. Roy[3 4]   Michael Carbin[1]

*Figure 1.* A diagram of instability analysis from step 0 (left) and step $k$ (right) when comparing networks using linear interpolation.

- Rewinding to iteration K, rather than init works much better
- Experimental analysis through the existence of linear connectivity
- Connectivity emerges early in training, but not at init (hard to find models that exhibit it)

# An interesting finding

# The value of values

**Deconstructing Lottery Tickets:**
**Zeros, Signs, and the Supermask**

**Hattie Zhou**
Uber
hattie@uber.com

**Janice Lan**
Uber AI
janlan@uber.com

**Rosanne Liu**
Uber AI
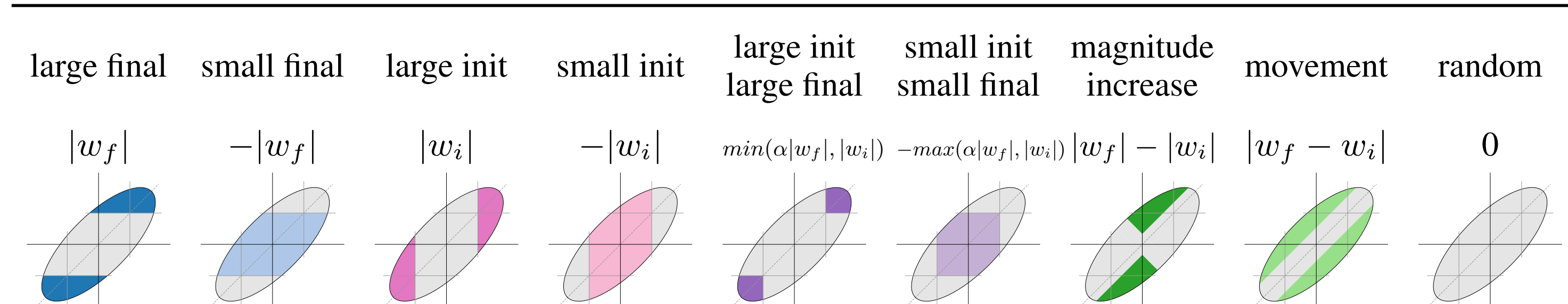rosanne@uber.com

**Jason Yosinski**
Uber AI
yosinski@uber.com

# Why do lottery tickets perform well?
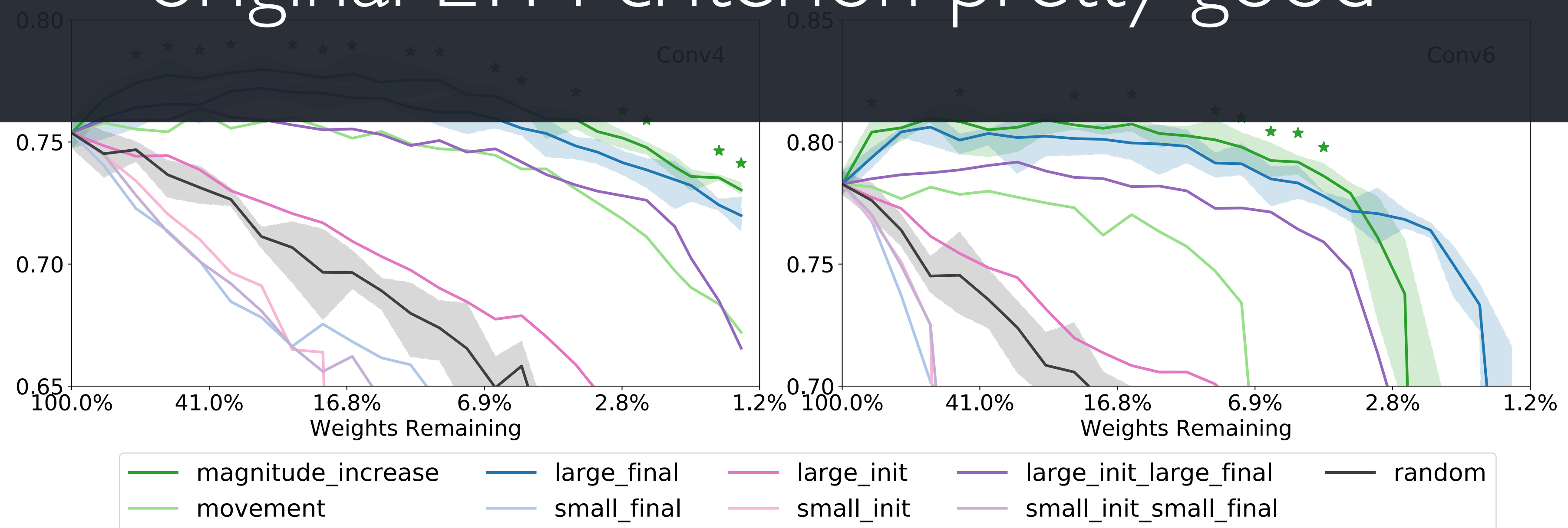
- *A Study on what allows LTs to be good*

Revisiting the IMP algorithm
1. Randomly initialize weights
2. pick "importance" metric M    **Mask criterion**
3. Train for small number of iterations
4. Prune bottom p% of according to M    **Mask-1 action**
5. Rollback top 100-p% non-zero weights    **Mask- 0 action**
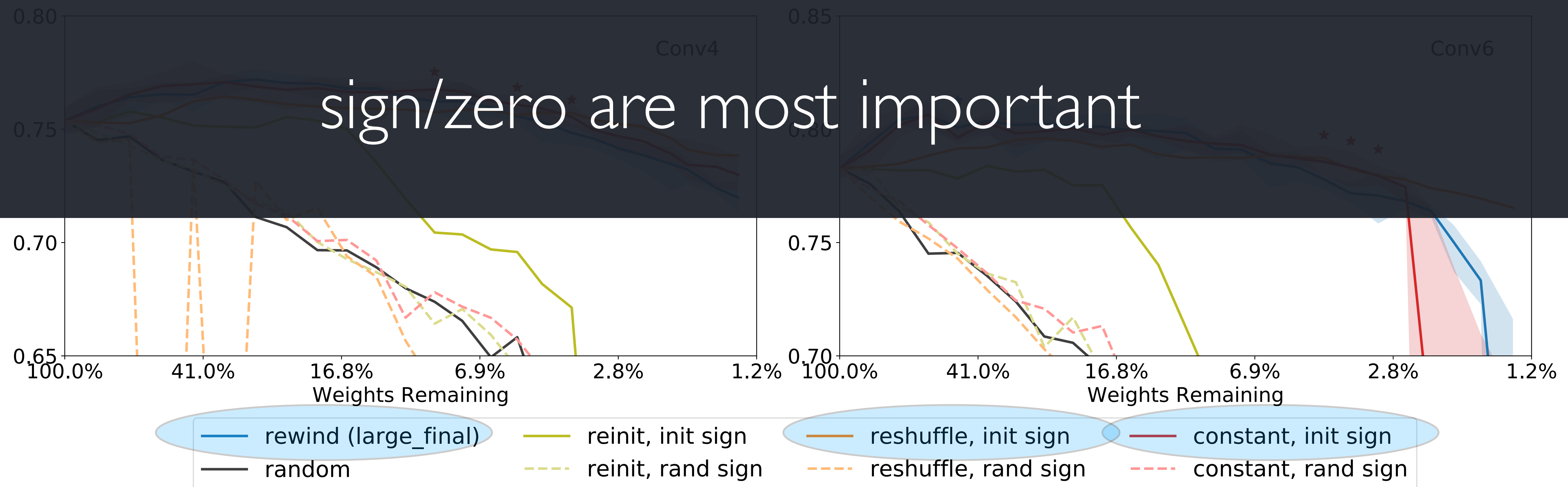6. re-train to full accuracy
7. (Sometimes) GOTO 3

# Mask Criteria

| large final | small final | large init | small init | large init large final | small init small final | magnitude increase | movement | random |
|---|---|---|---|---|---|---|---|---|
| $|w_f|$ | $-|w_f|$ | $|w_i|$ | $-|w_i|$ | $min(\alpha|w_f|, |w_i|)$ | $-max(\alpha|w_f|, |w_i|)$ | $|w_f| - |w_i|$ | $|w_f - w_i|$ | $0$ |

## original LTH criterion pretty good



Legend: magnitude_increase, movement, large_final, small_final, large_init, small_init, large_init_large_final, small_init_small_final, random

# What to do with surviving weights?

- *Reinitialize*
- *Shuffle original values*
- *Replace with sign(w)\*std(init. values)*



sign/zero are most important

# What to do with pruned weights?

- *freeze to init value*
- *set to zero*

## zero is special:
## learning the supermask similar to training

- *indication that pruning (without training) attains non-trivial test error*

# Pruning is all you need??

# What's Hidden in a Randomly Weighted Neural Network?

Vivek Ramanujan [*][†]     Mitchell Wortsman [*][‡]     Aniruddha Kembhavi [†][‡]

Ali Farhadi [‡]     Mohammad Rastegari [‡]

## Abstract

*Training a neural network is synonymous with learning the values of the weights. In contrast, we demonstrate that randomly weighted neural networks contain subnetworks which achieve impressive performance without ever modifying the weight values. Hidden in a randomly weighted Wide ResNet-50 [32] we find a subnetwork (with random weights) that is smaller than, but matches the performance of a ResNet-34 [9] trained on ImageNet [4]. Not only do these "untrained subnetworks" exist, but we provide an algorithm to effectively find them. We empirically show that as randomly weighted neural networks with fixed weights grow wider and deeper, an "untrained subnetwork" approaches a network with learned weights in accuracy. Our code and pretrained models are available at: https://github.com/allenai/hidden-networks.*
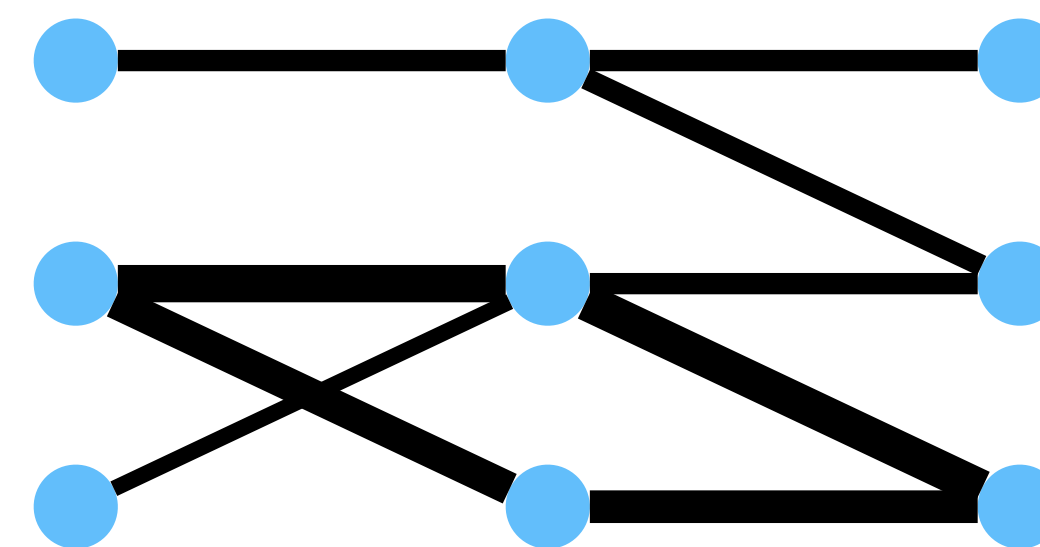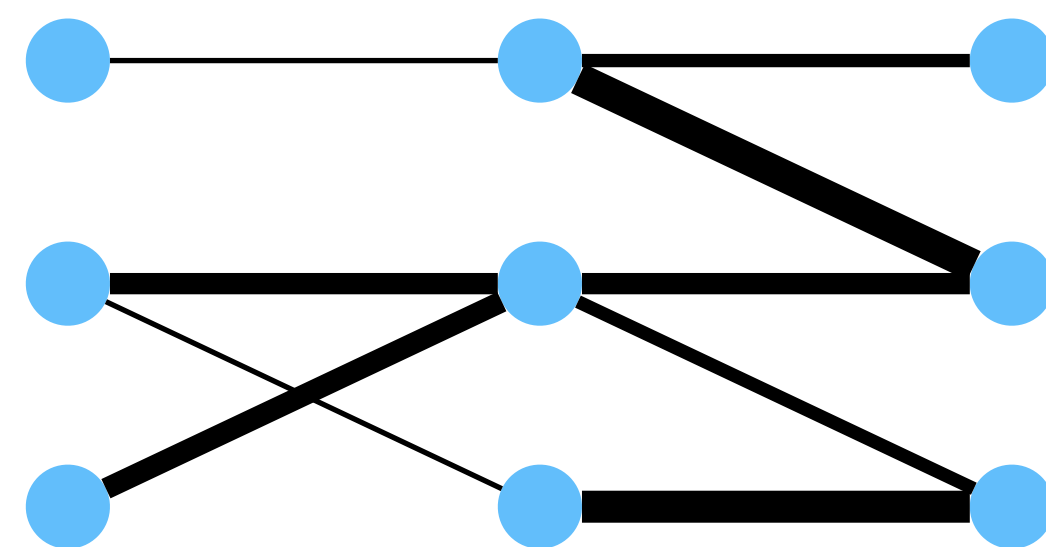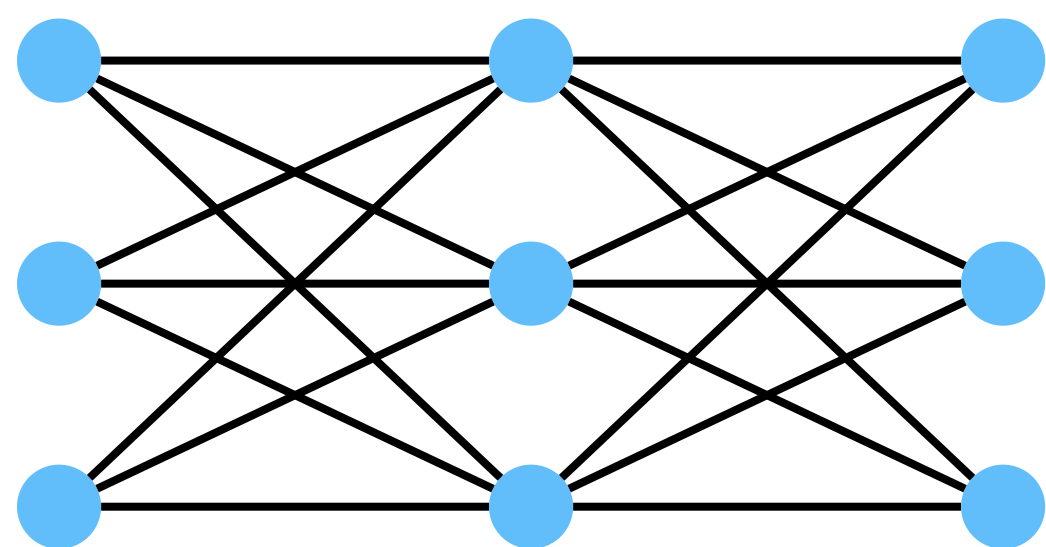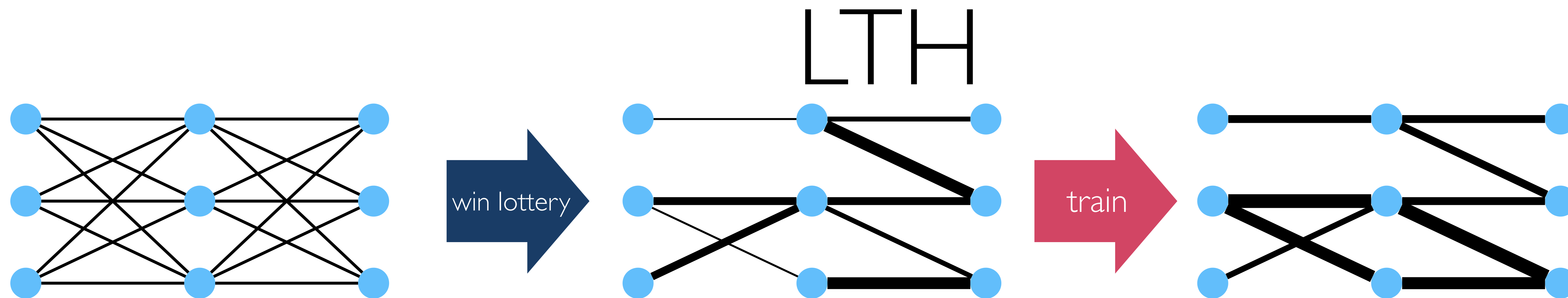
A neural network $\tau$ which achieves good performance

Randomly initialized neural network $N$
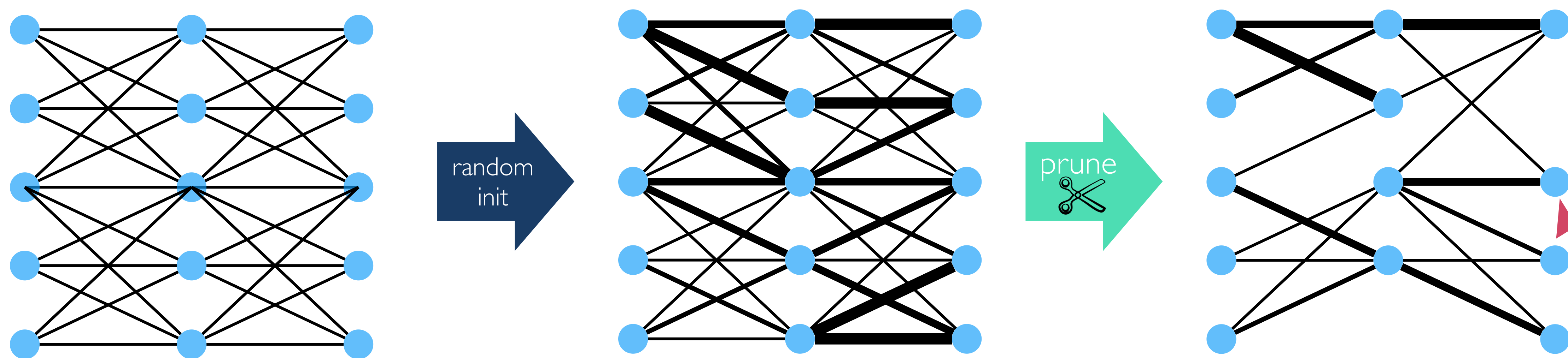
A subnetwork $\tau'$ of $N$

Figure 1. If a neural network with random weights (center) is sufficiently overparameterized, it will contain a subnetwork (right) that perform as well as a trained neural network (left) with the same number of parameters.

## 1. Introduction

# LTH

# LTH



win lottery

train

# Strong LTH: pruning is all you need

SOTA accuracy tickets simply reside within random NNs

random init

prune

≈

Woah, hold on…
You can get a high accuracy model…
WITHOUT SGD??

Q: Is the strong LTH universally true?

I.e., Can we always approximate a target NN
by pruning a larger random network?

well, if the larger net contains all possible weights..

# Conclusions & Open Problems

- "Trainable" sparse nets are desirable
- "early" *LTH = winning tickets exist at initialization*
- *later stage LTH = well, not quite, you have to train a bit first*
- *Finding LTs at init seems hard. Is it impossible though?*
- *Many extensions to BERT, Low-rank models, structured pruning*
- *Pruning is learning? WTFeta?*

Open Questions:
- Sparsity vs overparameterization
- Can we prune at initialization
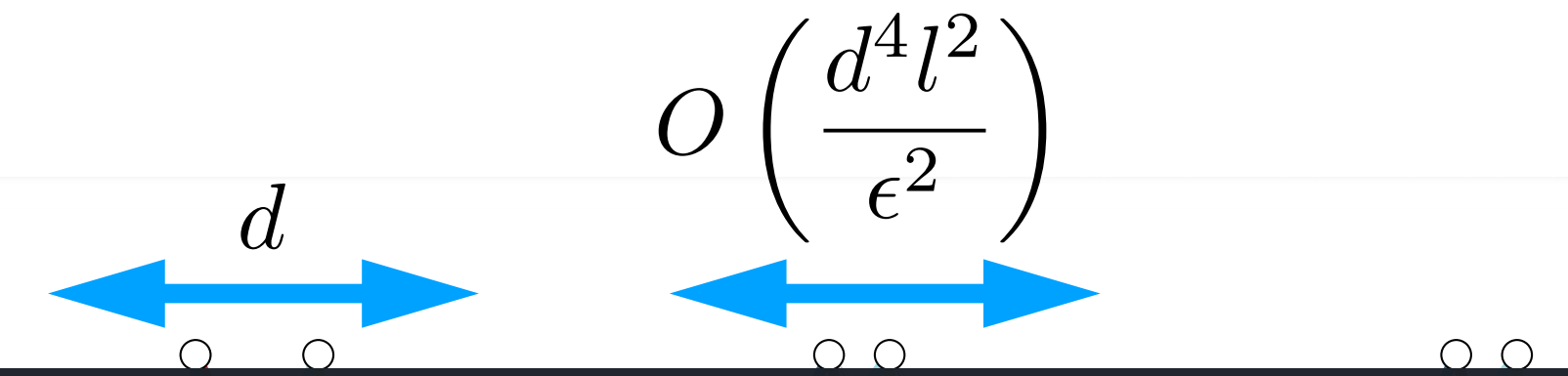- Where's the math??

# Part II: Theory
*(mostly existential results)*

# Do lottery tickets exist?

- Even in the absence of computational concerns, do LTs exist?

- If so, under what conditions?

- Provable poly-time algorithms?

# Malach et al. Proving the Strong LTH

$$O\left(\frac{d^4 l^2}{\epsilon^2}\right)$$

$d$

Proving the Lottery Ticket Hypothesis: Pruning is All You Need

Eran Malach[1], Gilad Yehudai[2], Shai Shalev-Shwartz[1], and Ohad Sh...

Weizmann Institute of Science

**Abstract**

The lottery ticket hypothesis (Frankle and Carbin, 2018), states that a random ... contains a small subnetwork such that, when trained in isolation, can compete with the performance of the original network. We prove an even stronger hypothesis (as was also conjectured in Ramanujan et al., 2019), showing that for every bounded distribution and every target network with bounded weights, a sufficiently over-parameterized neural network with random weights contains a subnetwork with roughly the same accuracy as the target network, without any further training.

A neural network $\tau$ which achieves good performance

Randomly initialized neural network $N$

A subnetwork $\tau'$ of $N$

Note:
This proves ANYTHING can be found in a larger net, e.g.,
vanilla LTs at init,
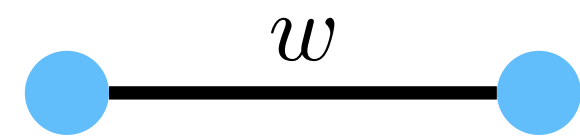later iteration LTs
"optimal" LTs

BUT… Ramanujan et al., prune a random WideResnet50
to approximate a Resnet 34
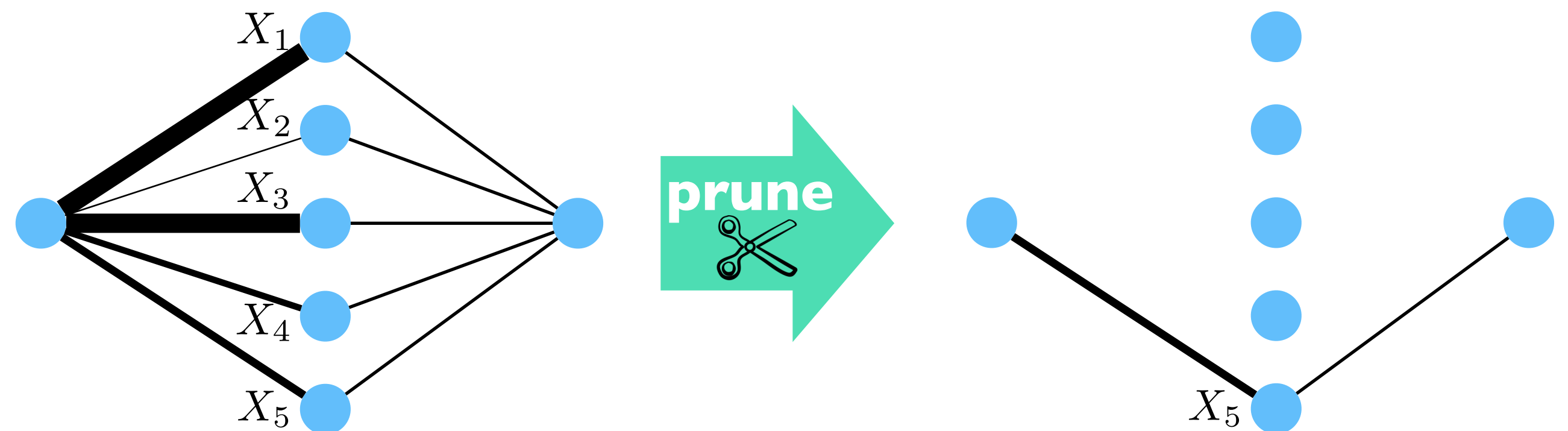
# Sketch of Malach et al.

Main idea:
If there are enough weights
one can approximately find the target NN

target weight

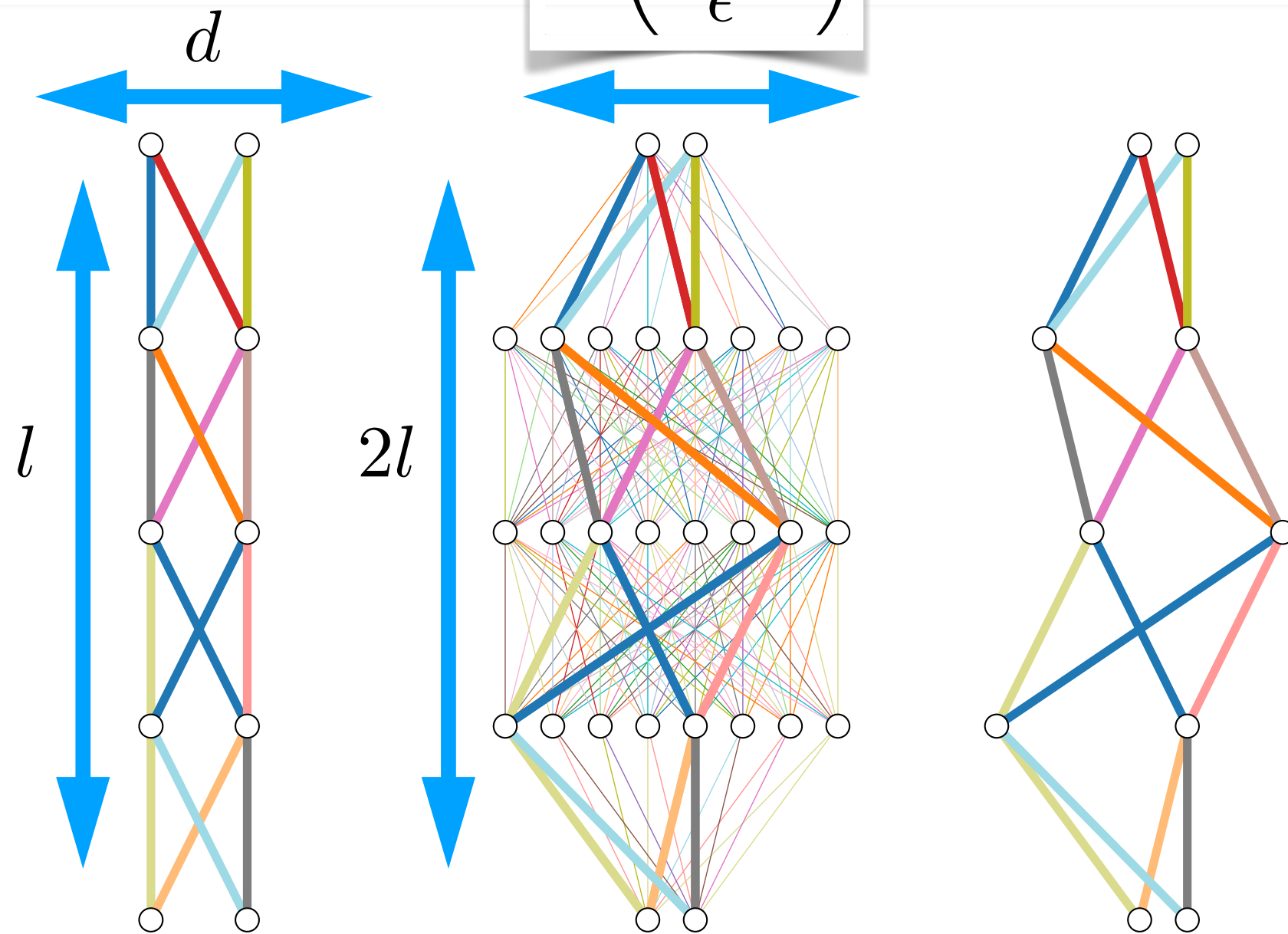pruning a highly over-parameterized network

$w$



$X_i \sim$ Uniform R.V

**Lemma**: if we draw $1/\varepsilon$ random weights,
one will be $\varepsilon$-close to target with constant. probability

The general theorem is a more involved extension of this idea

So we can't improve on this?

**poly**$(1/\varepsilon)$ dependence unavoidable, if you prune down to one weight

# Our work: Exponentially tighter Strong LTH



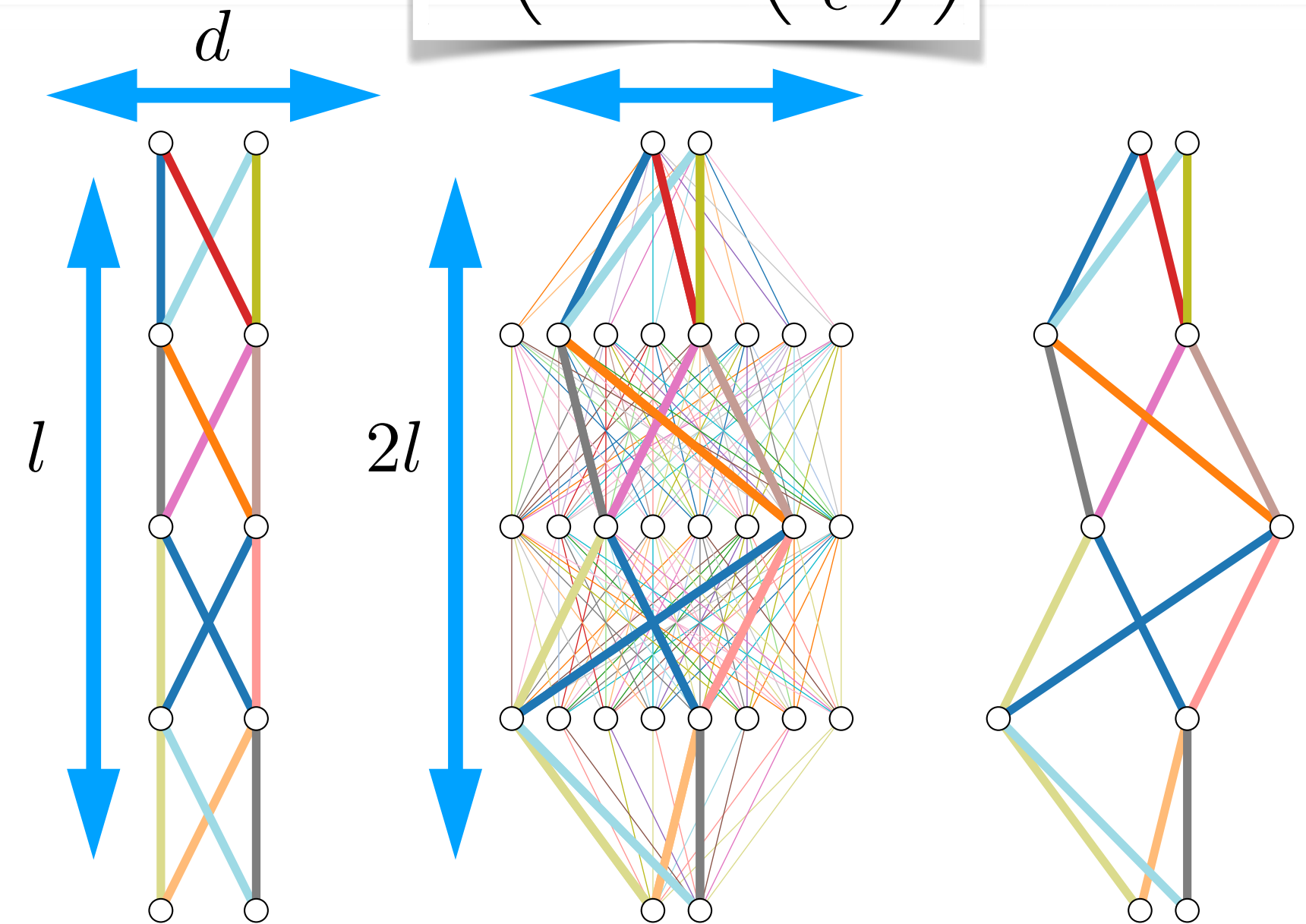Malach et al.

$$O\left(\frac{d^4 l^2}{\epsilon^2}\right)$$

$d$

$l$

$2l$

A neural network $\tau$ which achieves good performance

Randomly initialized neural network $N$

A subnetwork $\tau'$ of $N$

winning a single weight lottery

Our work

$$O\left(d \cdot \log\left(\frac{dl}{\epsilon}\right)\right)$$

$d$

$l$

$2l$

A neural network $\tau$ which achieves good performance

Randomly initialized neural network $N$

A subnetwork $\tau'$ of $N$

winning a subset span lottery

# The Subset Span approach

# Malach theorem = pruning ε-nets



$X_i \sim$ Uniform R.V

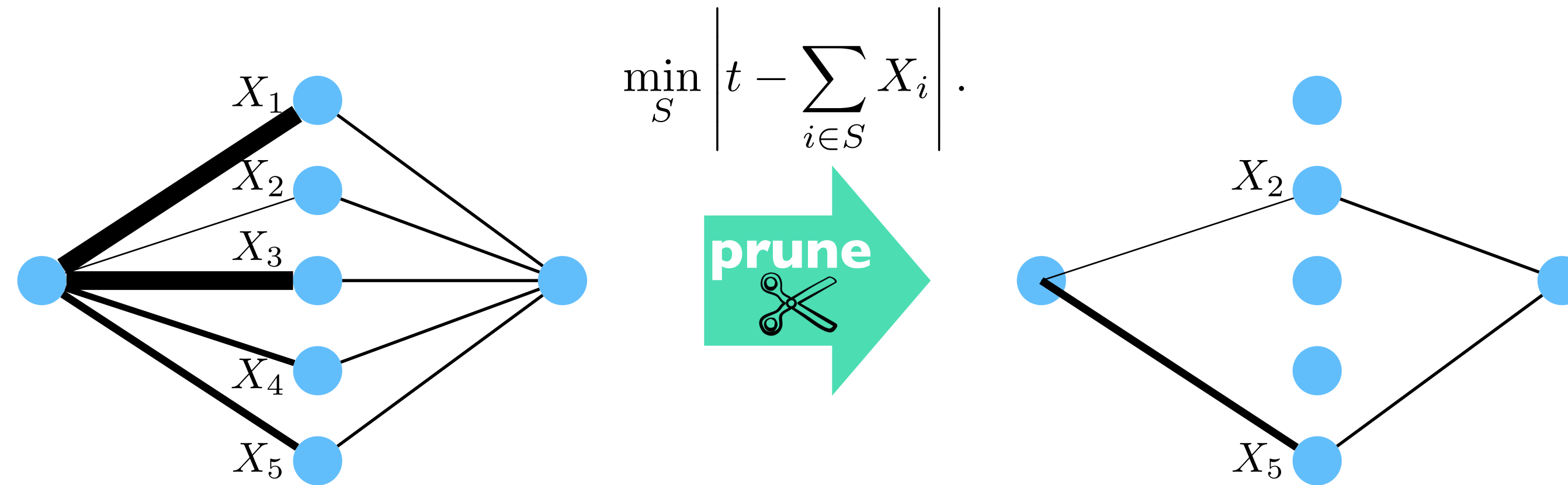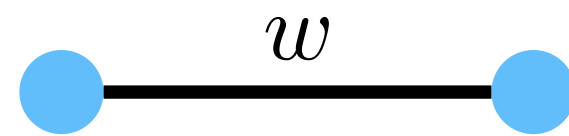| generating enough so that any number falls ε-close |
|---|

| pruning = finding closest point to ε-net |
|---|

**poly**(1/ε) dependence unavoidable

what if we combine subsets of those random weights?

# SUBSET Span

target weight

$w$

$$\min_S \left| t - \sum_{i \in S} X_i \right|.$$

$X_1$ $X_2$ $X_3$ $X_4$ $X_5$

**prune**

$X_2$ $X_5$

pruning = finding best subset sum to approximate target

Note: this is like a "batch" version of single parameter pruning

**Q:** how many RVs do I need for an ε-approximation?

# [Lueker 1998]

## Exponentially Small Bounds on the Expected Optimum of the Partition and Subset Sum Problems*

*George S. Lueker*

Department of Information and Computer Science
University of California, Irvine
Irvine, CA  92697-3425

**Theorem 2.4.** *Let $X_1, X_2, \ldots, X_n$ be i.i.d. uniform over $[-1, 1]$, and let $0 < \eta < \frac{1}{2}$. Suppose that $n/2 \geq C \ln \eta^{-1}$. Then, except with probability bounded by*
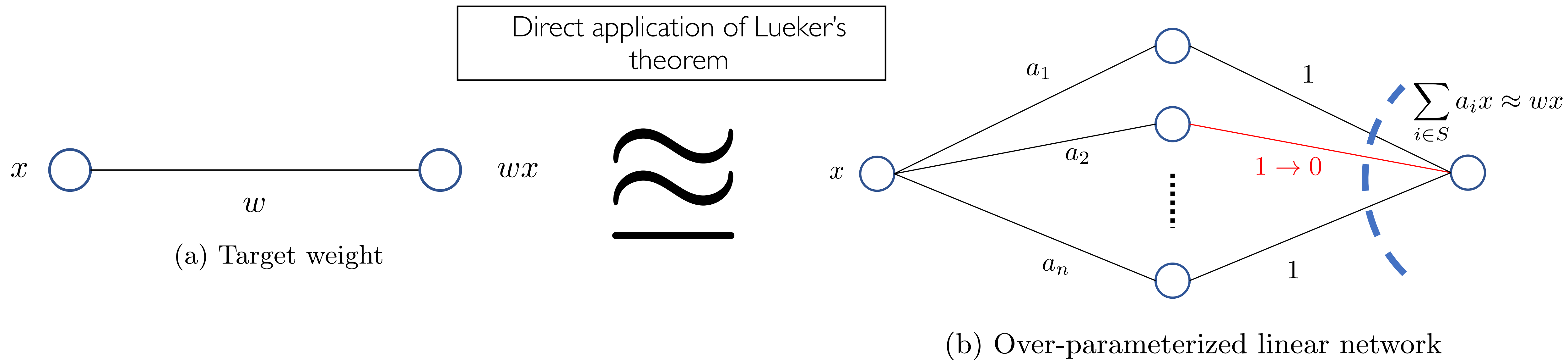
$$\exp\left(-\frac{\left(n/2 - C \ln \eta^{-1}\right)^2}{2n}\right),$$

*all values in $[-\frac{1}{2}, \frac{1}{2}]$ have admissible $2\eta$-approximations.*
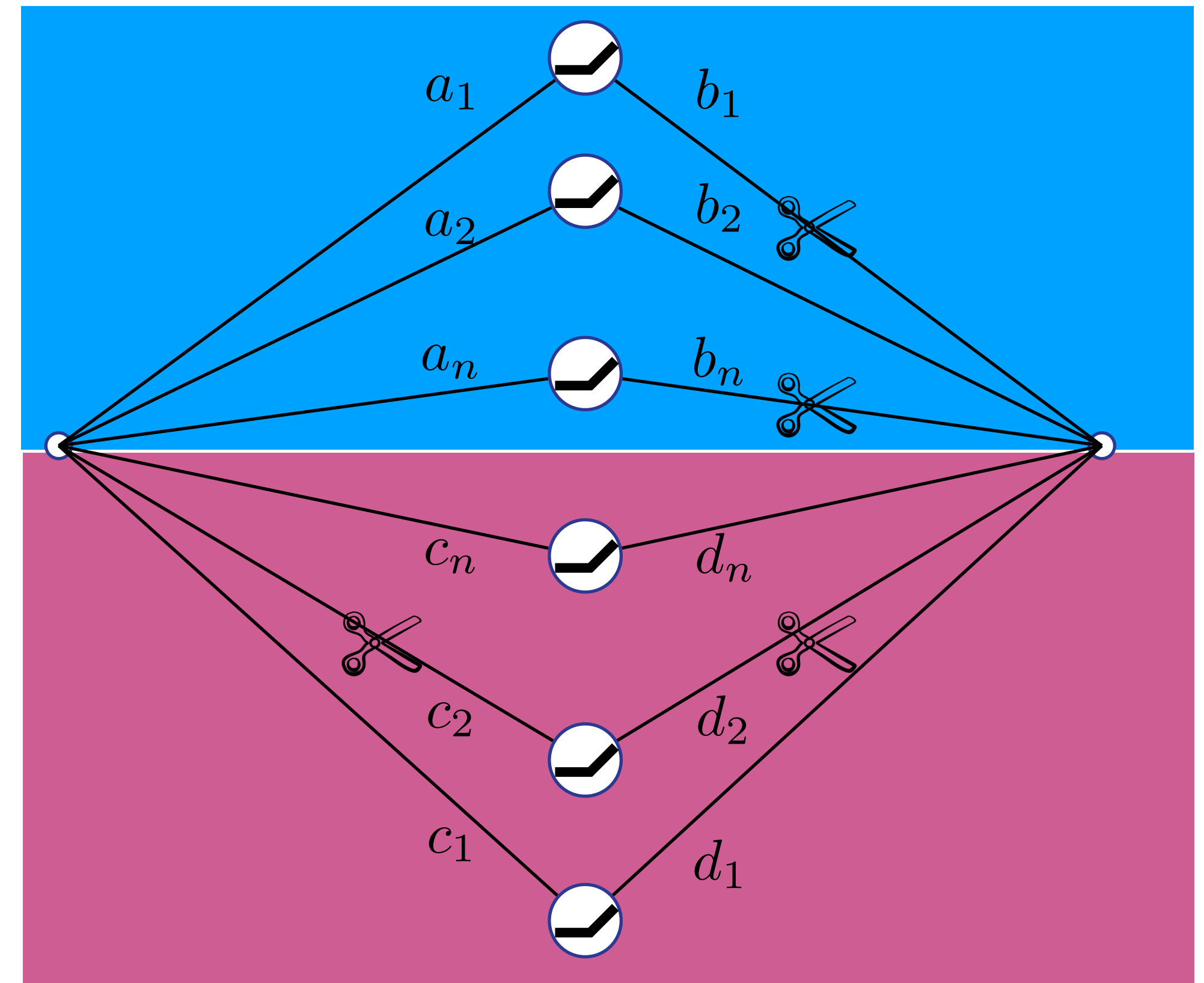
The Subset span is very expressive:
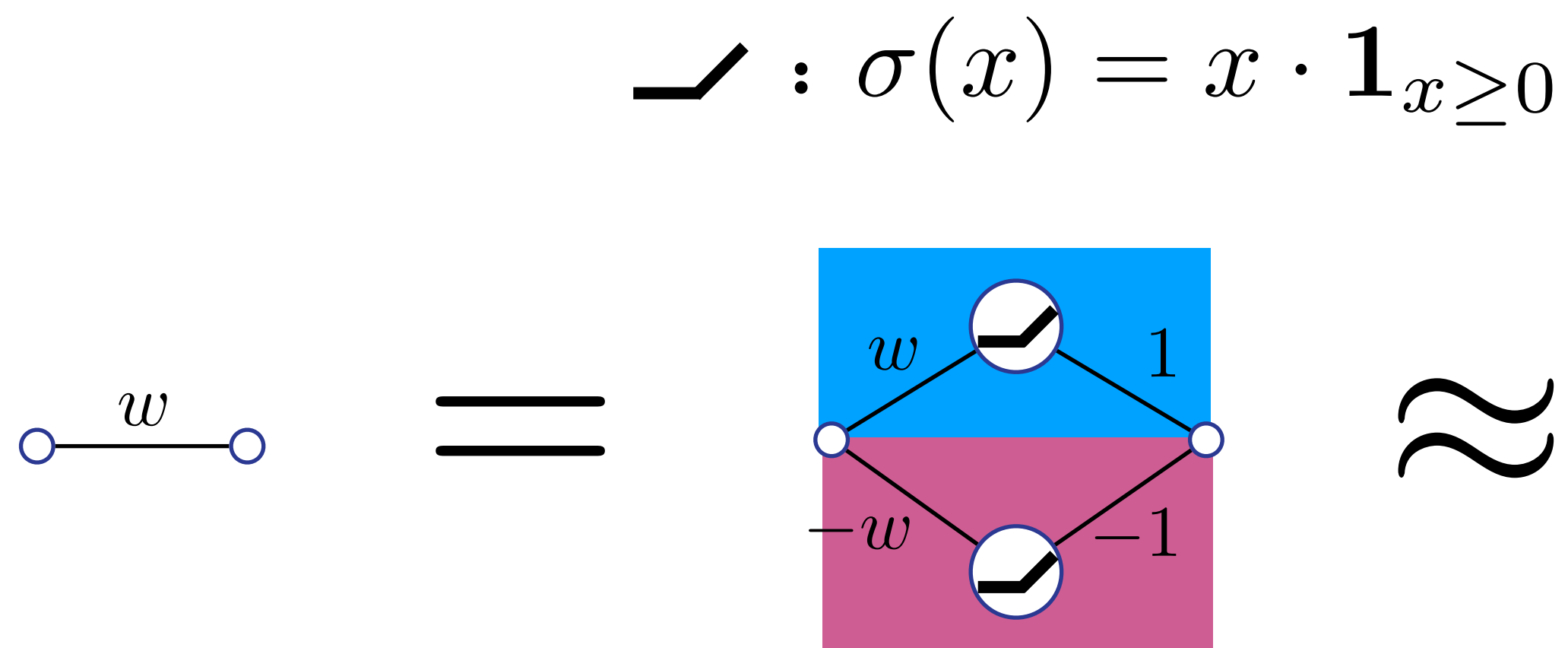
Every number in [-1,1], can be approximated by taking a subset of $log(1/\varepsilon)$ RVs

# How to approximate a single weight

Direct application of Lueker's theorem



(a) Target weight

$$\approx$$

$$\sum_{i \in S} a_i x \approx wx$$

$1 \to 0$

(b) Over-parameterized linear network

How do we transform the linear net to a ReLu?
Constraint: Weights have to be uniform and iid.

# How to approximate a single weight

$$\diagup \; : \; \sigma(x) = x \cdot \mathbf{1}_{x \geq 0}$$
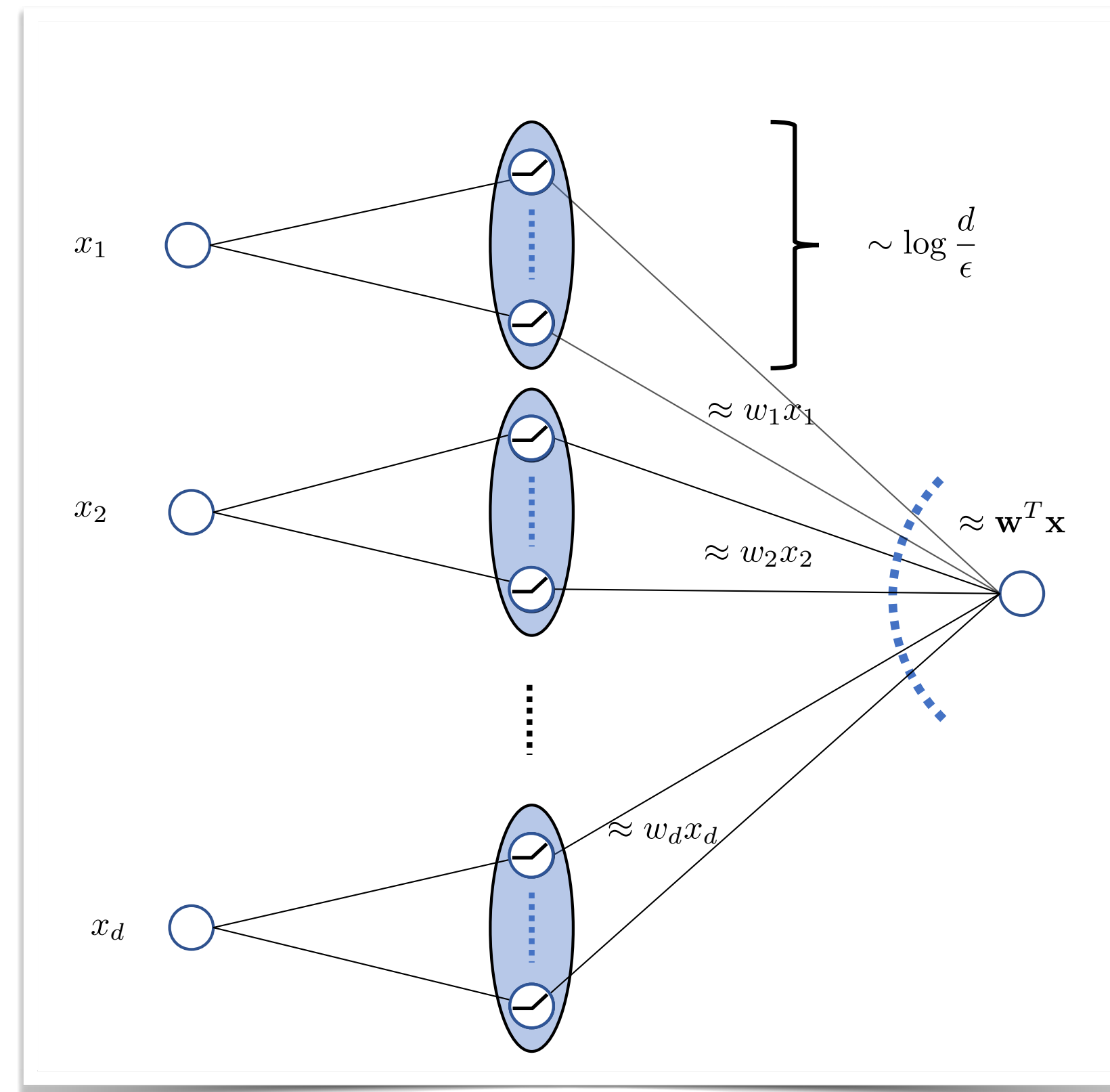


Lueker's theorem still holds if the distribution of a*b  contain Uniform[-1,1]

# from weights to neurons to networks

$$w \cdot x$$

$$\mathbf{w}^T\mathbf{x}$$

$$\mathbf{W}^T\mathbf{x}$$



$$\sum_{i \in S_1} b_i \sigma(a_i x) \approx \sigma(wx)$$

$$\sum_{i \in S_2} d_i \sigma(c_i x) \approx -\sigma(-wx)$$

$$\sim \log \frac{d}{\epsilon}$$

$$\approx w_1 x_1$$

$$\approx w_2 x_2$$

$$\approx \mathbf{w}^T\mathbf{x}$$

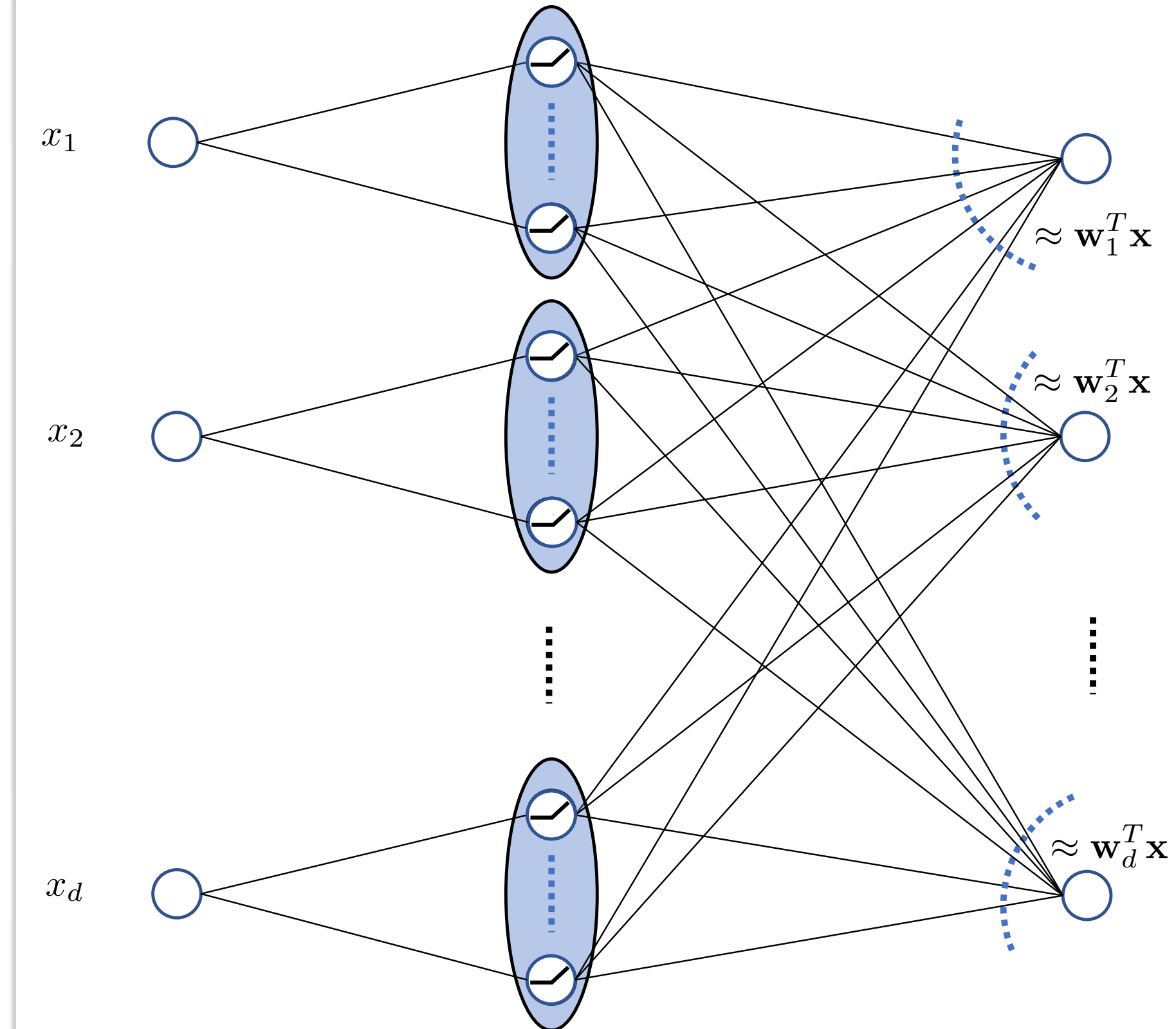$$\approx w_d x_d$$

$$\approx \mathbf{w}_1^T\mathbf{x}$$

$$\approx \mathbf{w}_2^T\mathbf{x}$$

$$\approx \mathbf{w}_d^T\mathbf{x}$$

All that is left: operator norm bounds for each approximation layer so that

$$\min_{\mathbf{S}_i \in \{0,1\}^{d_i \times d_{i-1}}} \ \sup_{\|\mathbf{x}\| \leq 1} \ \|f(\mathbf{x}) - (\mathbf{S}_{2l} \odot \mathbf{M}_{2l})\sigma((\mathbf{S}_{2l-1} \odot \mathbf{M}_{2l-1}) \ldots \sigma((\mathbf{S}_1 \odot \mathbf{M}_1)\mathbf{x}))\| < \epsilon.$$

# Lower bound via Packing

**Theorem 2.** *(informal) There exists a 2-layer neural network with width d which cannot be approximated to error within $\epsilon$ by pruning a randomly initialized 2-layer network, unless the random network has width at least $\Omega(d\log(1/\epsilon))$.*
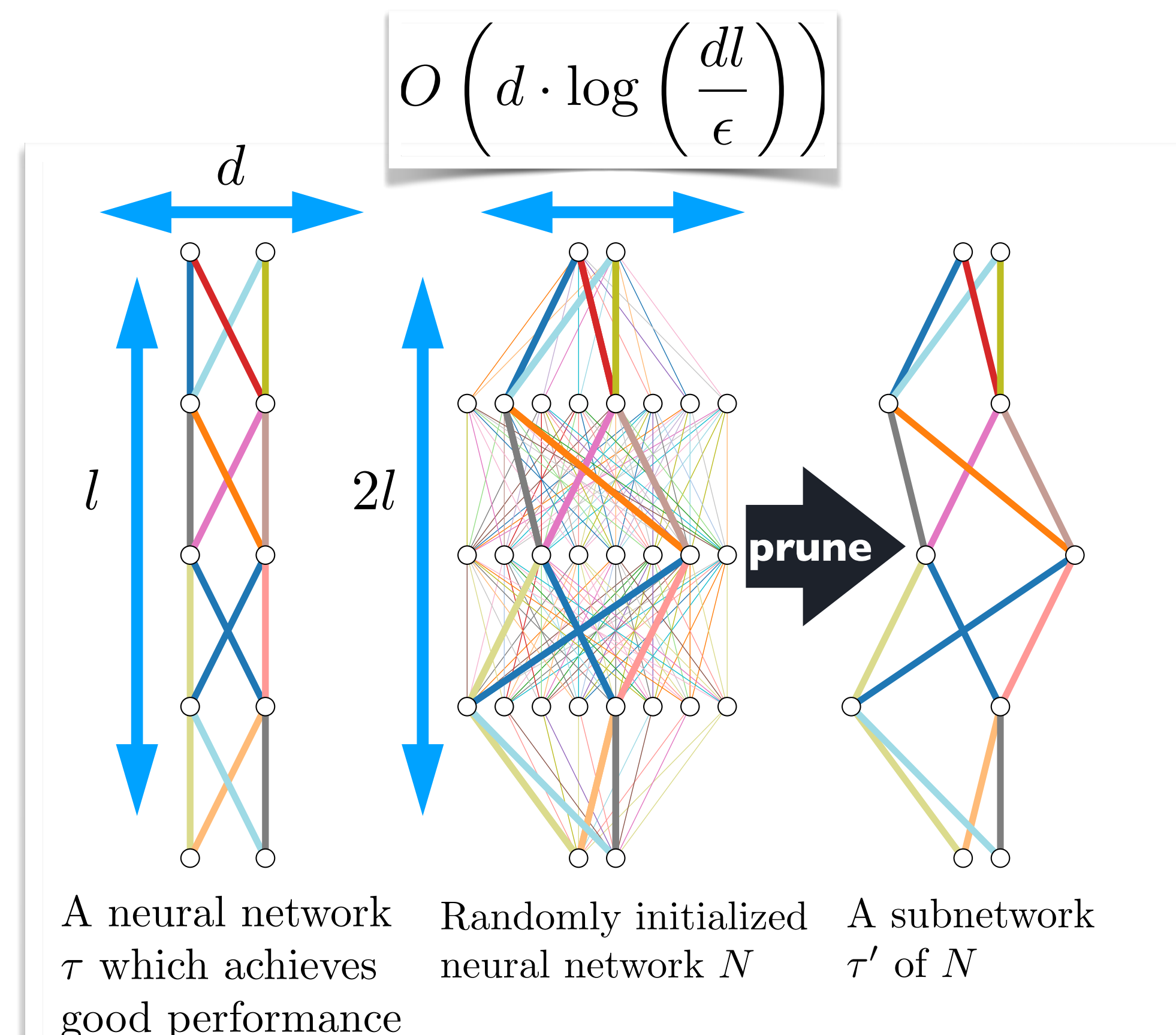
Proof idea:
How many $\varepsilon$-separated linear functions can we pack
in a large pruned matrix?

# Learning ⊂ Pruning

Any neural network be approximated by pruning
a logarithmically overparameterized network of random* weights

$$O\left(d \cdot \log\left(\frac{dl}{\epsilon}\right)\right)$$



$d$

$l$  $2l$

**prune**

A neural network
$\tau$ which achieves
good performance

Randomly initialized
neural network $N$
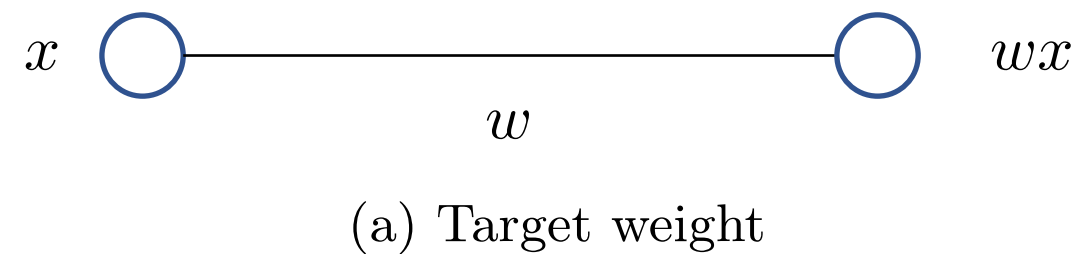
A subnetwork
$\tau'$ of $N$

Note:
this is an existential result. Although
our proof is algorithmic, we do not
propose a new pruning algorithm

one experiment

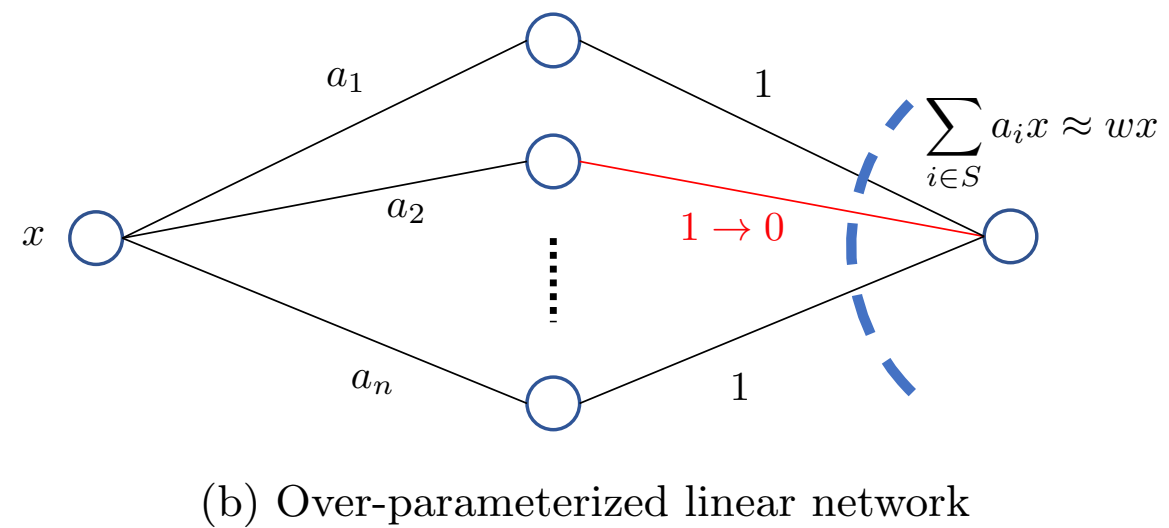# Type of Overparam Matters a lot!!
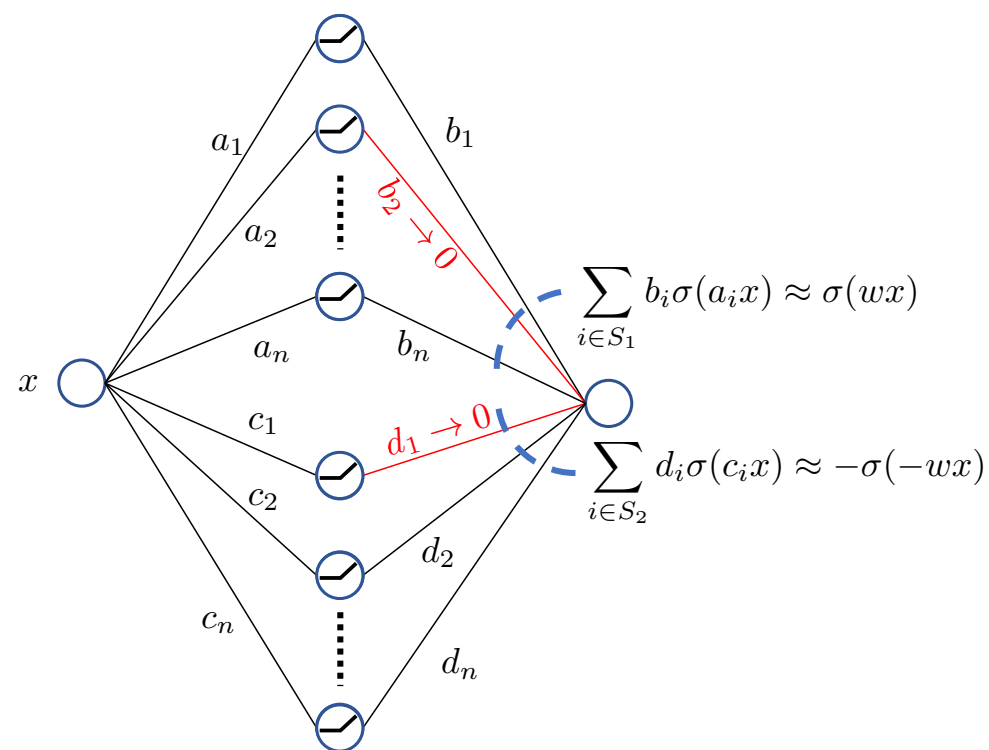
- Comparison for pruning wider nets
  - FC ReLU nets



(a) Target weight

  - Nets with linear "diamond" structure



(b) Over-parameterized linear network

  - Nets with ReLU "diamond" structure



(c) Over-parameterized ReLU network



Pruning via the Ramanujan et al. algorithm

Accuracy vs Number of Parameters

Legend:
- Our Structure, Linear
- Our Structure, ReLU
- Wide Network
- Baseline

(c) LeNet5

# Conclusions & Open Problems

- A DlogD random net contains ALL networks of size D!
- Vanilla LTs exist! So do Perfect LTs!
- The IMP's problem is not existence, but algorithmic.
- One can learn by pruning

Open Question:
- Can we fix IMP?
- Prune + train existential results?
- Can pruning be faster than training? (better for hardware?)
- Network architectures amenable to pruning
- Towards a "no-backprop" training framework

# Reading List

Han, S., Mao, H. and Dally, W.J., 2015. Deep compression: Compressing deep neural networks with pruning, trained quantization and huffman coding. arXiv preprint arXiv:1510.00149.

Frankle, J. and Carbin, M., 2018. The lottery ticket hypothesis: Finding sparse, trainable neural networks. ICLR 2019

Malach, E., Yehudai, G., Shalev-Schwartz, S. and Shamir, O., 2020, November. Proving the lottery ticket hypothesis: Pruning is all you need. In International Conference on Machine Learning (pp. 6682-6691). PMLR.

Frankle, J., Dziugaite, G.K., Roy, D. and Carbin, M., 2020, November. Linear mode connectivity and the lottery ticket hypothesis. In International Conference on Machine Learning (pp. 3259-3269). PMLR.

Zhou, H., Lan, J., Liu, R. and Yosinski, J., 2019. Deconstructing lottery tickets: Zeros, signs, and the supermask. Advances in neural information processing systems, 32.

Liu, Z., Sun, M., Zhou, T., Huang, G. and Darrell, T., 2018. Rethinking the value of network pruning. arXiv preprint arXiv:1810.05270.

Gale, T., Elsen, E. and Hooker, S., 2019. The state of sparsity in deep neural networks. arXiv preprint arXiv:1902.09574.

Ramanujan, V., Wortsman, M., Kembhavi, A., Farhadi, A. and Rastegari, M., 2020. What's hidden in a randomly weighted neural network?. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (pp. 11893-11902).

Pensia, A., Rajput, S., Nagle, A., Vishwakarma, H. and Papailiopoulos, D., 2020. Optimal lottery tickets via subset sum: Logarithmic over-parameterization is sufficient. Advances in Neural Information Processing Systems, 33, pp.2599-2610.