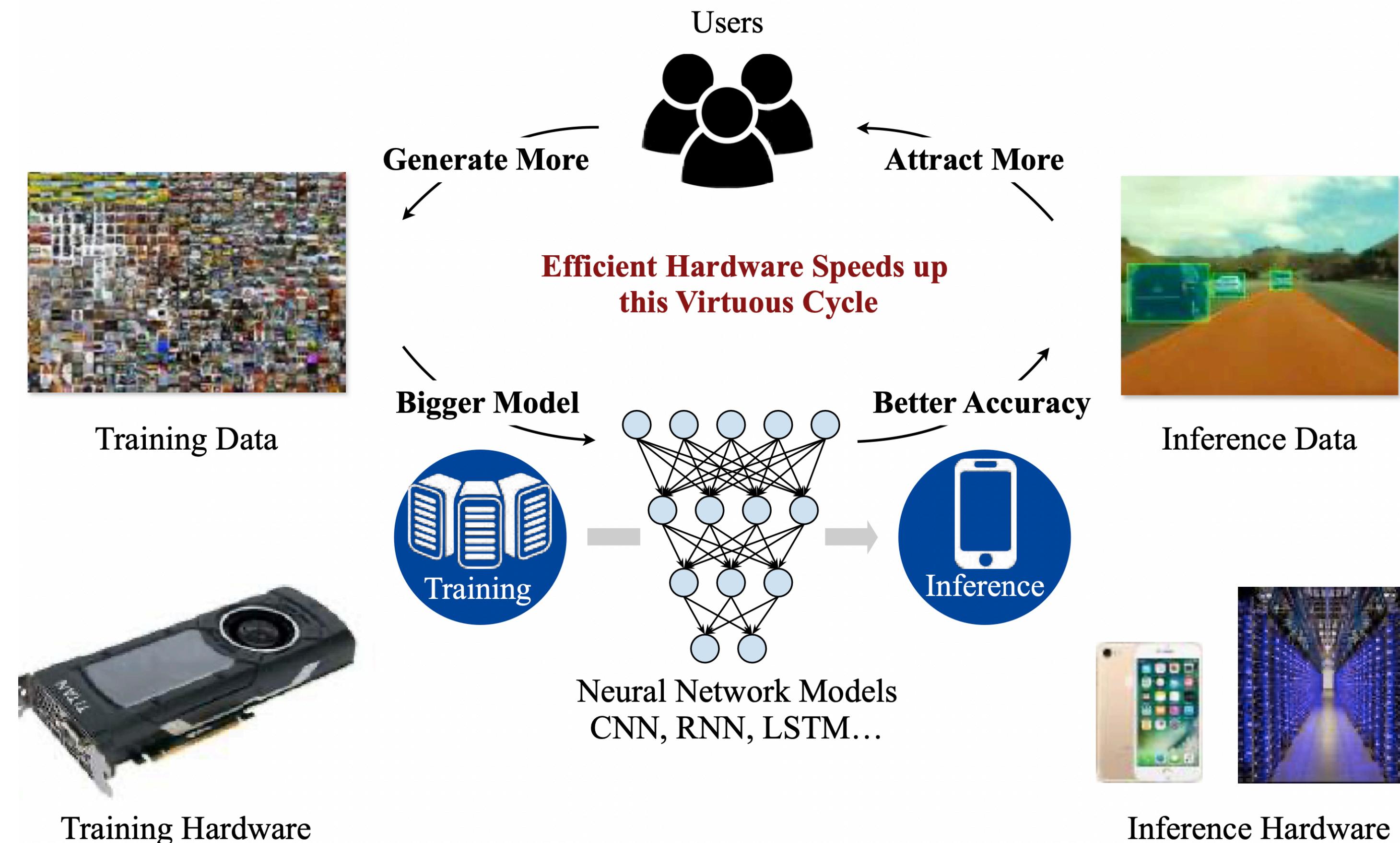


# **Deep Compression: Compressing Deep Neural networks with Pruning, Trained Quantization and Huffman Coding**

**Song Han, Huizi Mao, and William J. Dally**

**Presenter: Muyang Li**

# Hardware Plays an Important Role in Modern Deep Learning



# A Challenge for Modern Deep Learning



- We are solving more complicated AI problems with larger datasets, which **requires more computation**.
- However, Moore's Law is slowing down; the amount of computation per unit cost is **no longer increasing** at its historic rate.

# We Do Need AI on Edge Devices

Low  
Latency

Low  
Cost

User  
Privacy



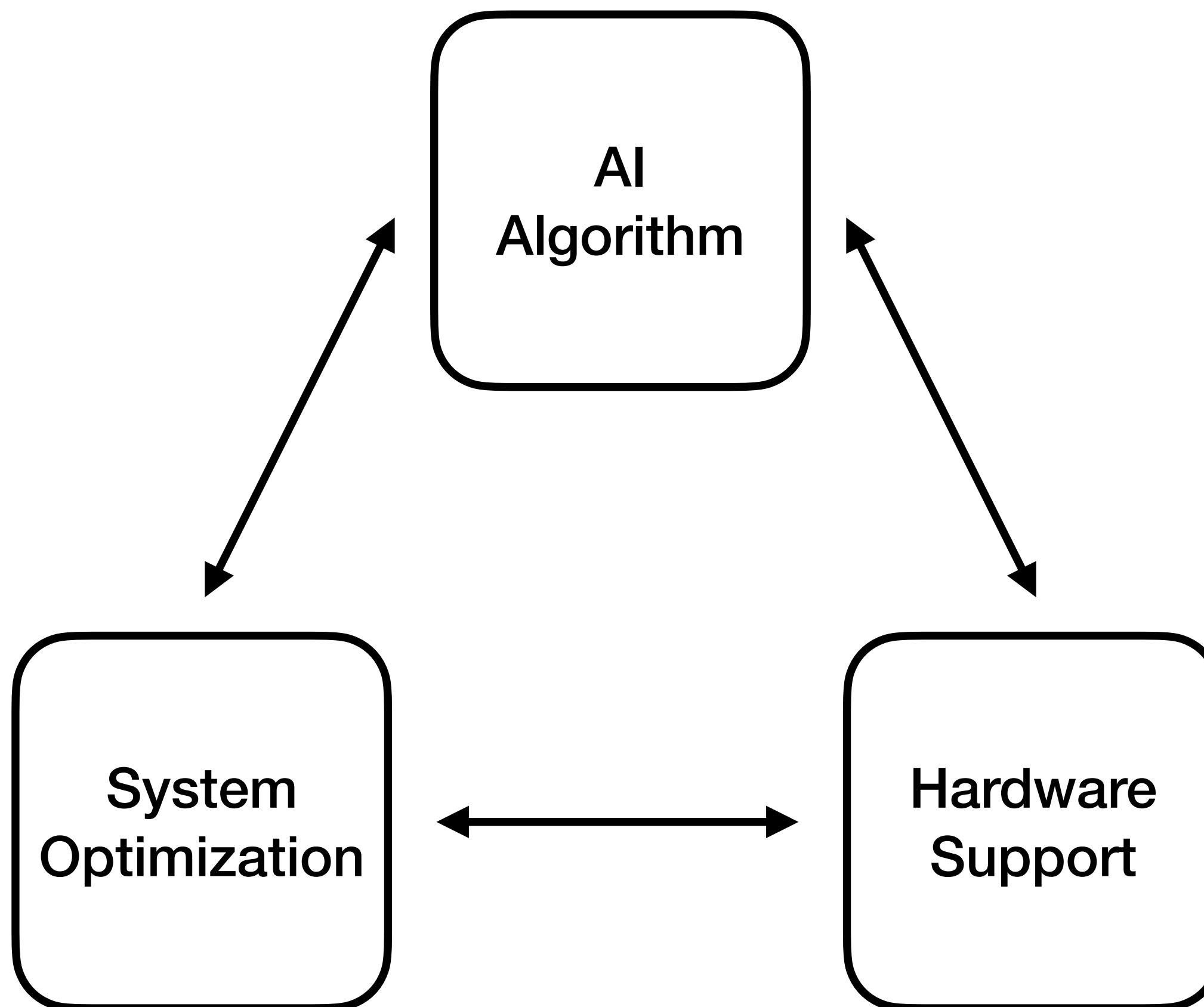
Edge devices only have limited resources!!!



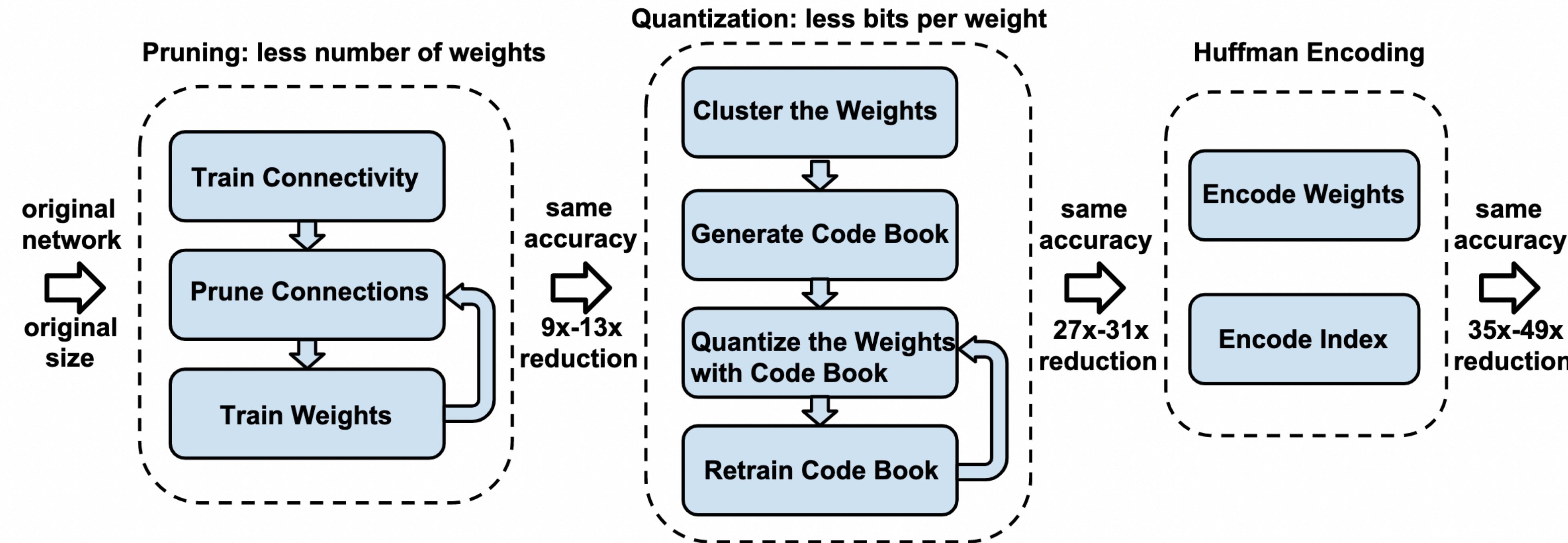
Power-  
Hungry

AI is good, but AI is slow!

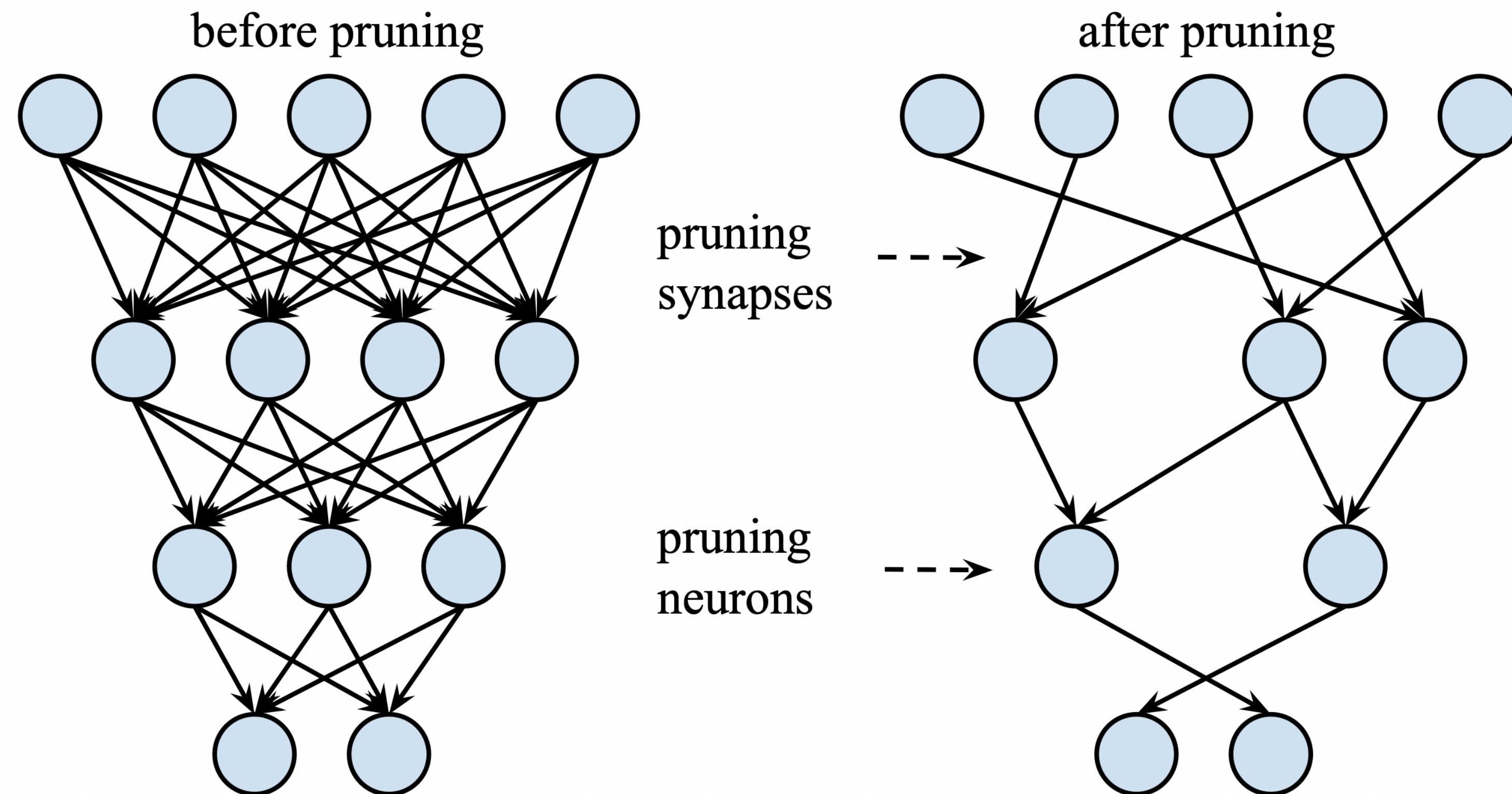
# Ways to Accelerate AI



# Deep Compression Overview



# Pruning



# Pruning: Sparse Representation

[https://en.wikipedia.org/wiki/Sparse matrix](https://en.wikipedia.org/wiki/Sparse_matrix) and Deep Compression, Han et al., 2015, <https://arxiv.org/pdf/1510.00149.pdf>

# Pruning: Sparse Representation

$$\begin{pmatrix} 10 & 20 & 0 & 0 & 0 & 0 \\ 0 & 30 & 0 & 0 & 0 & 40 \\ 0 & 0 & 50 & 60 & 70 & 0 \\ 0 & 0 & 0 & 0 & 0 & 80 \end{pmatrix}$$

Original Matrix

[https://en.wikipedia.org/wiki/Sparse\\_matrix](https://en.wikipedia.org/wiki/Sparse_matrix) and Deep Compression, Han et al., 2015, <https://arxiv.org/pdf/1510.00149.pdf>

# Pruning: Sparse Representation

$$\begin{pmatrix} 10 & 20 & 0 & 0 & 0 & 0 \\ 0 & 30 & 0 & 0 & 0 & 40 \\ 0 & 0 & 50 & 60 & 70 & 0 \\ 0 & 0 & 0 & 0 & 0 & 80 \end{pmatrix}$$

Original Matrix

$$V = [10 \ 20 \ 30 \ 40 \ 50 \ 60 \ 70 \ 80]$$

$$COL = [0 \ 1 \ 1 \ 5 \ 2 \ 3 \ 4 \ 5]$$

$$ROW = [0 \ 2 \ 4 \ 7 \ 8]$$

CSR Representation ( $2a + n + 1$  numbers)

# Pruning: Sparse Representation

$$\begin{pmatrix} 10 & 20 & 0 & 0 & 0 & 0 \\ 0 & 30 & 0 & 0 & 0 & 40 \\ 0 & 0 & 50 & 60 & 70 & 0 \\ 0 & 0 & 0 & 0 & 0 & 80 \end{pmatrix}$$

Original Matrix

$$V = [10 \ 20 \ 30 \ 40 \ 50 \ 60 \ 70 \ 80]$$

$$COL = [0 \ 1 \ 1 \ 5 \ 2 \ 3 \ 4 \ 5]$$

$$ROW = [0 \ 2 \ 4 \ 7 \ 8]$$

CSR Representation ( $2a + n + 1$  numbers)

$$V = [10 \ 20 \ 30 \ 40 \ 50 \ 60 \ 70 \ 80]$$

$$COL = [0 \ 1 \ 1 \ 4 \ 2 \ 1 \ 1 \ 5]$$

$$ROW = [0 \ 2 \ 4 \ 7 \ 8]$$

CSR Representation with Difference

# Pruning: Sparse Representation

Span Exceeds  $8=2^3$

A table illustrating a sparse representation. The columns are indexed from 0 to 15. The first row contains 'idx' values, the second row contains 'diff' values, and the third row contains 'value' (non-zero) values. The 'value' row has yellow background for indices 1, 4, and 15, and a grey background for index 12. A red arrow points to the value at index 12, which is 0, with the label 'Filler Zero' below it. A text annotation 'Span Exceeds 8=2^3' with arrows spans across the first 8 columns.

idx	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
diff		1			3								8			3
value		3.4			0.9								0			1.7

$$V = [10 \ 20 \ 30 \ 40 \ 50 \ 60 \ 70 \ 80]$$

$$COL = [0 \ 1 \ 1 \ 4 \ 2 \ 1 \ 1 \ 5]$$

$$ROW = [0 \ 2 \ 4 \ 7 \ 8]$$

CSR Representation with Difference

# Pruning: Sparse Representation

$$\begin{pmatrix} 10 & 20 & 0 & 0 & 0 & 0 \\ 0 & 30 & 0 & 0 & 0 & 40 \\ 0 & 0 & 50 & 60 & 70 & 0 \\ 0 & 0 & 0 & 0 & 0 & 80 \end{pmatrix}$$

Original Matrix

$$V = [10 \ 20 \ 30 \ 40 \ 50 \ 60 \ 70 \ 80]$$

$$COL = [0 \ 1 \ 1 \ 5 \ 2 \ 3 \ 4 \ 5]$$

$$ROW = [0 \ 2 \ 4 \ 7 \ 8]$$

CSR Representation ( $2a + n + 1$  numbers)

$$V = [10 \ 20 \ 30 \ 40 \ 50 \ 60 \ 70 \ 80]$$

$$COL = [0 \ 1 \ 1 \ 4 \ 2 \ 1 \ 1 \ 5]$$

$$ROW = [0 \ 2 \ 4 \ 7 \ 8]$$

CSR Representation with Difference

# Pruning: Sparse Representation

$$\begin{pmatrix} 10 & 20 & 0 & 0 & 0 & 0 \\ 0 & 30 & 0 & 0 & 0 & 40 \\ 0 & 0 & 50 & 60 & 70 & 0 \\ 0 & 0 & 0 & 0 & 0 & 80 \end{pmatrix}$$

Original Matrix

$$V = [10 \ 20 \ 30 \ 40 \ 50 \ 60 \ 70 \ 80]$$

$$COL = [0 \ 1 \ 1 \ 5 \ 2 \ 3 \ 4 \ 5]$$

$$ROW = [0 \ 2 \ 4 \ 7 \ 8]$$

CSR Representation ( $2a + n + 1$  numbers)

$$V = [10 \ 20 \ 30 \ 40 \ 50 \ 60 \ 70 \ 80]$$

$$COL = [0 \ 1 \ 1 \ 4 \ 2 \ 1 \ 1 \ 5]$$

$$ROW = [0 \ 2 \ 4 \ 7 \ 8]$$

CSR Representation with Difference

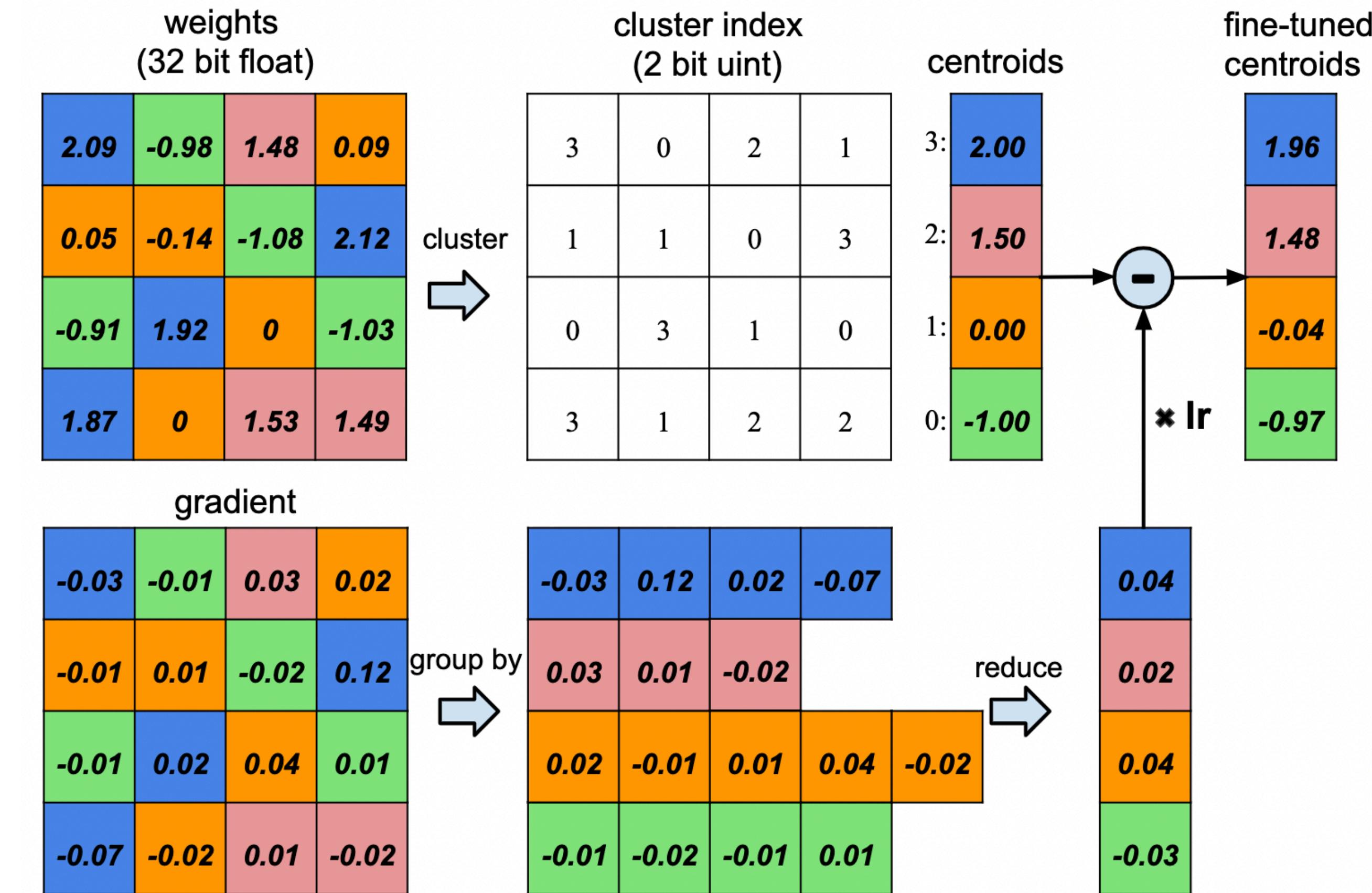
$$V = [10 \ 20 \ 30 \ 0 \ 40 \ 50 \ 60 \ 70 \ 0 \ 80]$$

$$COL = [0 \ 1 \ 1 \ 3 \ 1 \ 2 \ 1 \ 1 \ 3 \ 2]$$

$$ROW = [0 \ 2 \ 5 \ 8 \ 10]$$

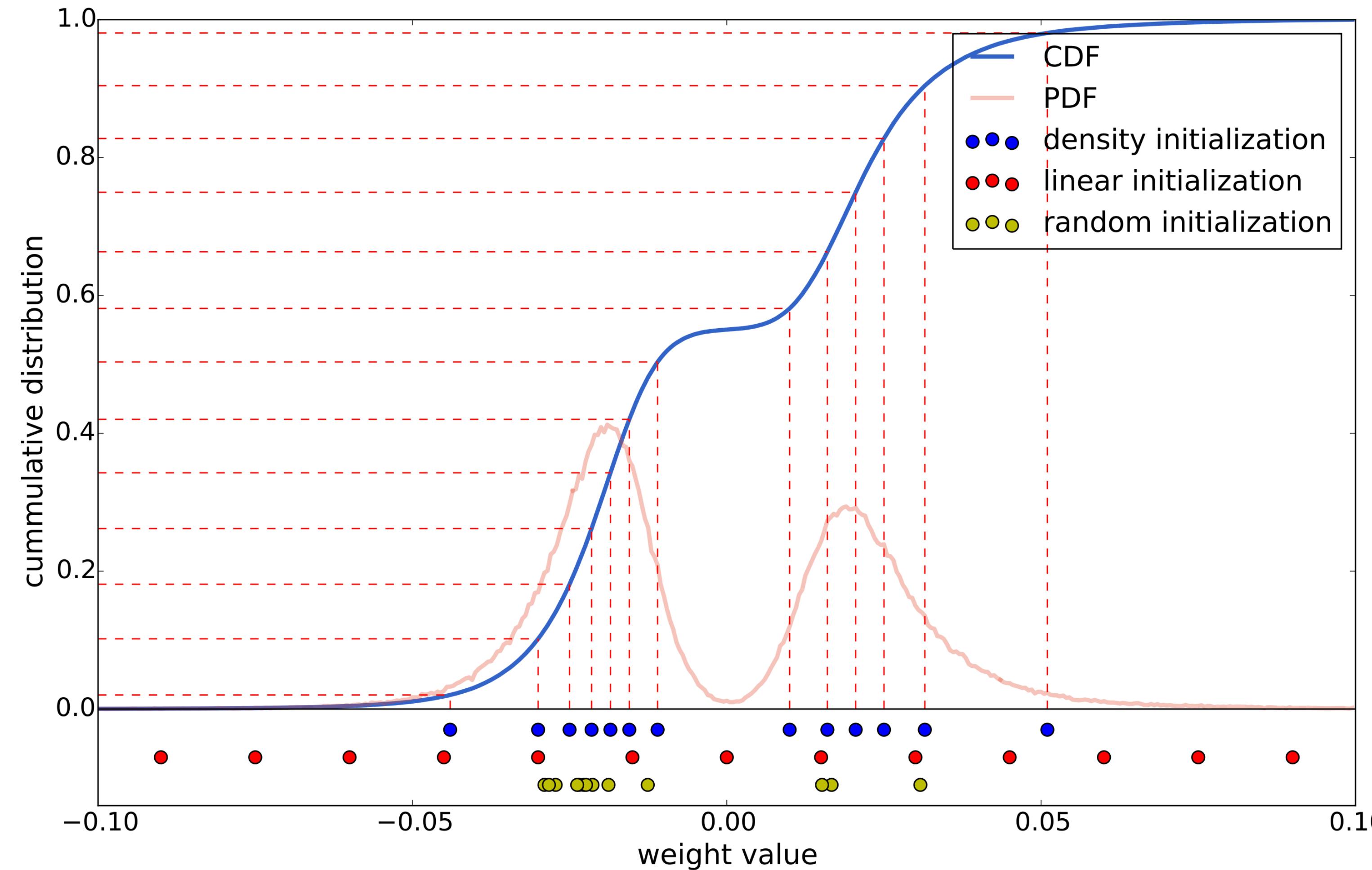
CSR Representation with Difference (encoded in 2 bits)

# Quantization

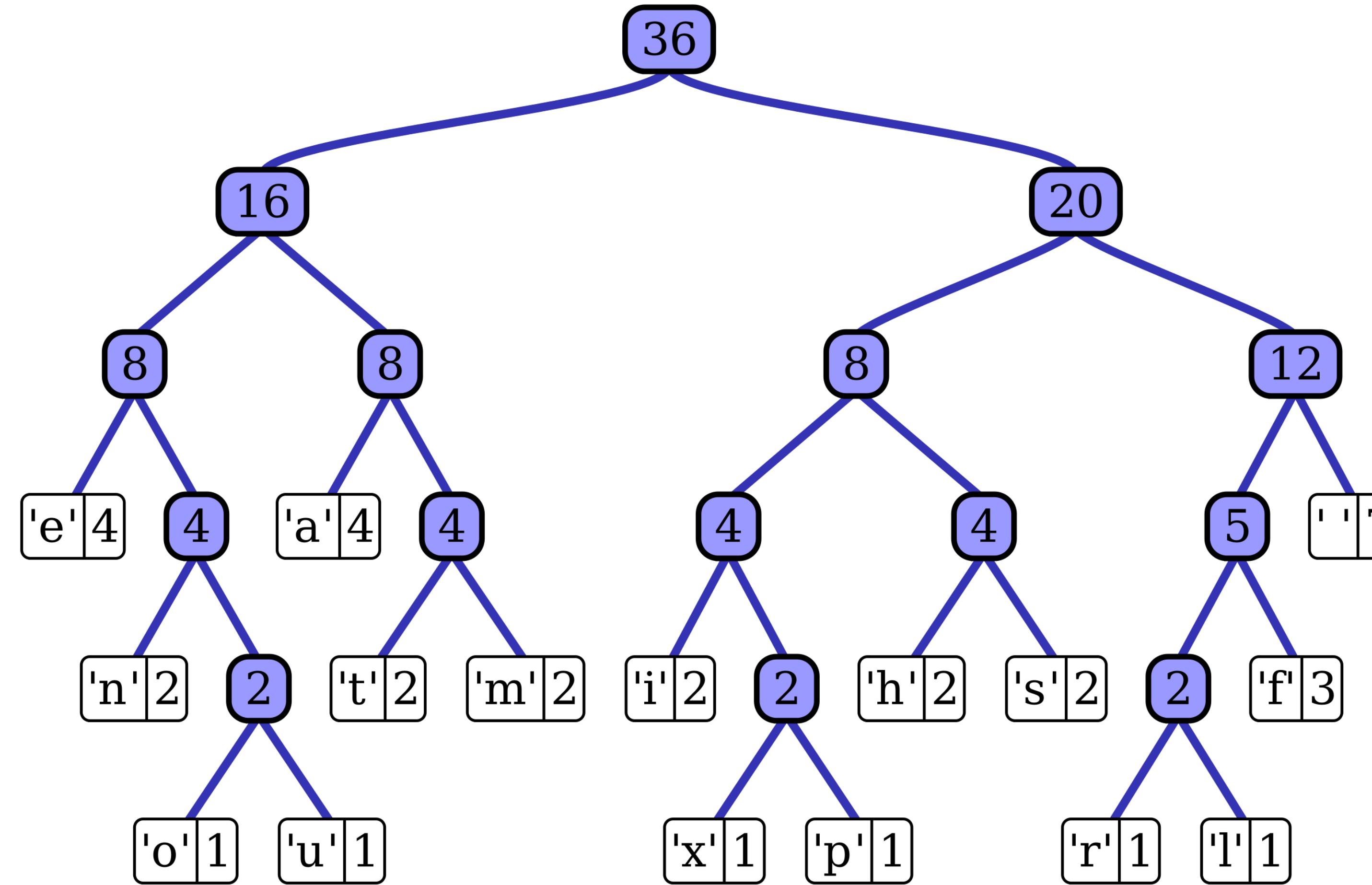


$$\text{Compression Ratio} \frac{nb}{n \log_2 k + kb}$$

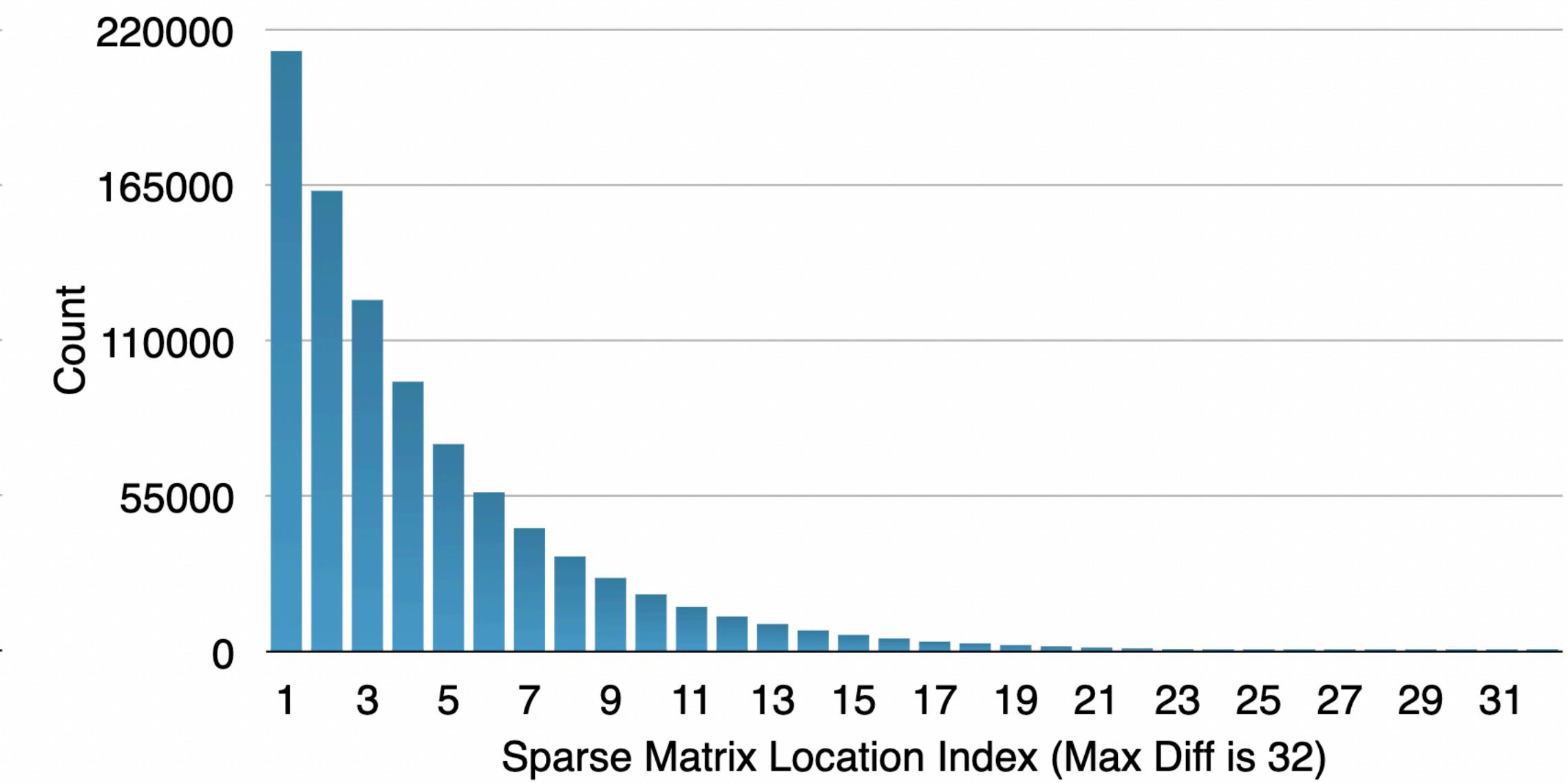
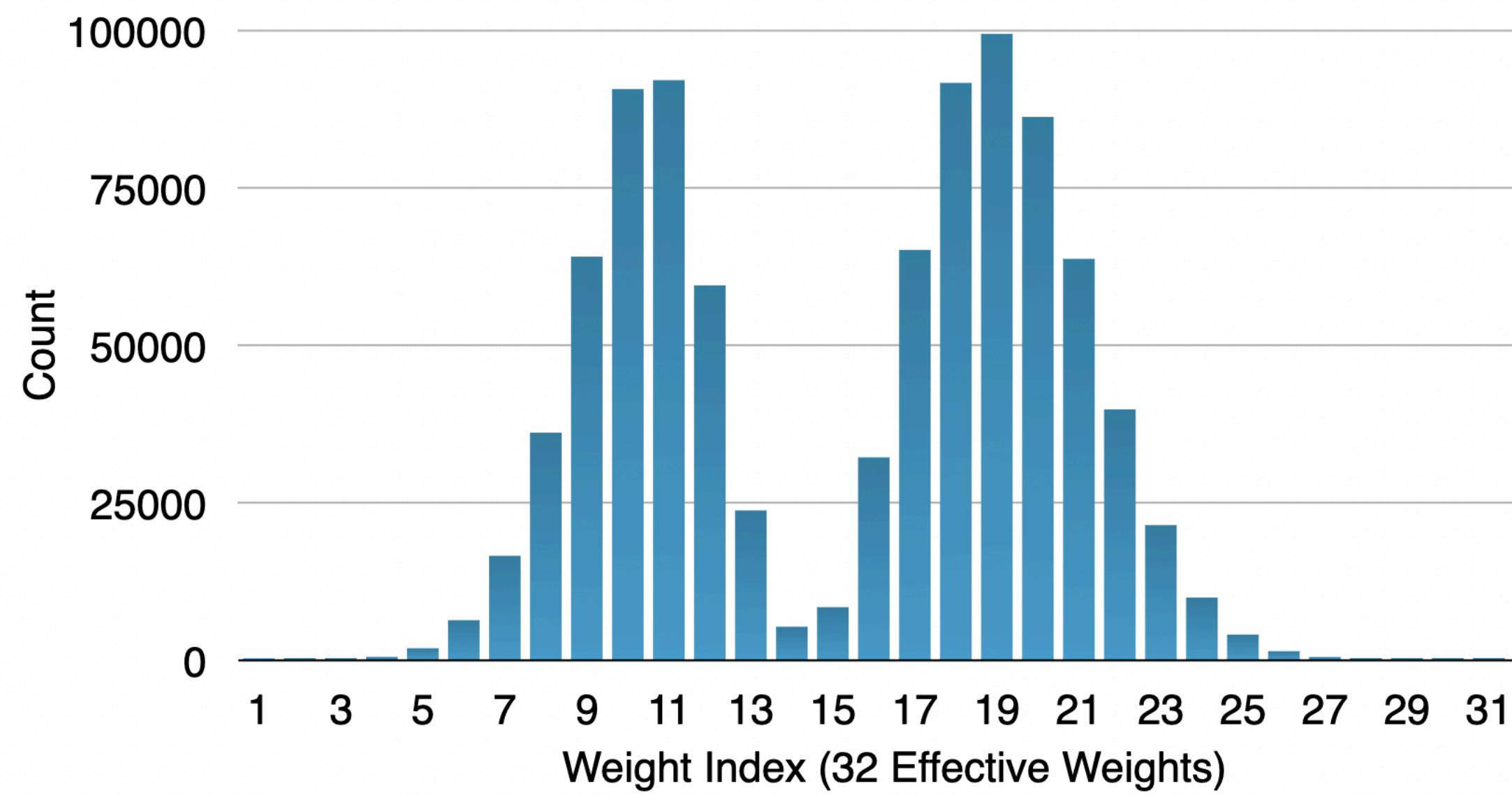
# Quantization: Initialization of the Shared Weights



# Huffman Coding



# Huffman Coding



# Results

Table 1: The compression pipeline can save  $35\times$  to  $49\times$  parameter storage with no loss of accuracy.

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

# Results

Network	Top-1 Error	Top-5 Error	Parameters	Compress Rate
Baseline Caffemodel ( <b>BVLC</b> )	42.78%	19.73%	240MB	1×
Fastfood-32-AD (Yang et al., 2014)	41.93%	-	131MB	2×
Fastfood-16-AD (Yang et al., 2014)	42.90%	-	64MB	3.7×
Collins & Kohli (Collins & Kohli, 2014)	44.40%	-	61MB	4×
SVD (Denton et al., 2014)	44.02%	20.56%	47.6MB	5×
Pruning (Han et al., 2015)	42.77%	19.67%	27MB	9×
Pruning+Quantization	42.78%	19.70%	8.9MB	27×
<b>Pruning+Quantization+Huffman</b>	<b>42.78%</b>	<b>19.70%</b>	<b>6.9MB</b>	<b>35×</b>

# Results

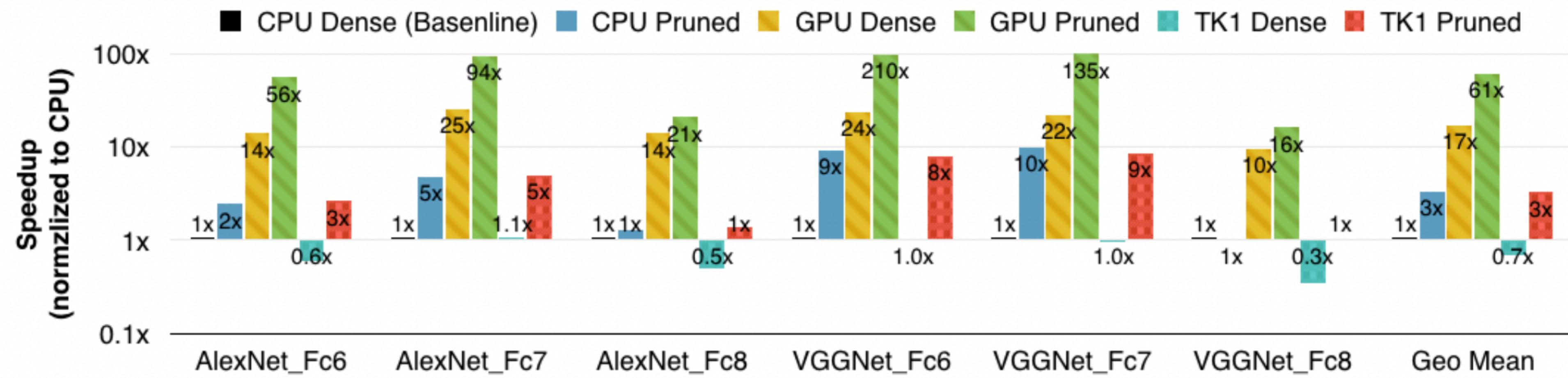


Figure 9: Compared with the original network, pruned network layer achieved  $3\times$  speedup on CPU,  $3.5\times$  on GPU and  $4.2\times$  on mobile GPU on average. Batch size = 1 targeting real time processing. Performance number normalized to CPU.

# Results

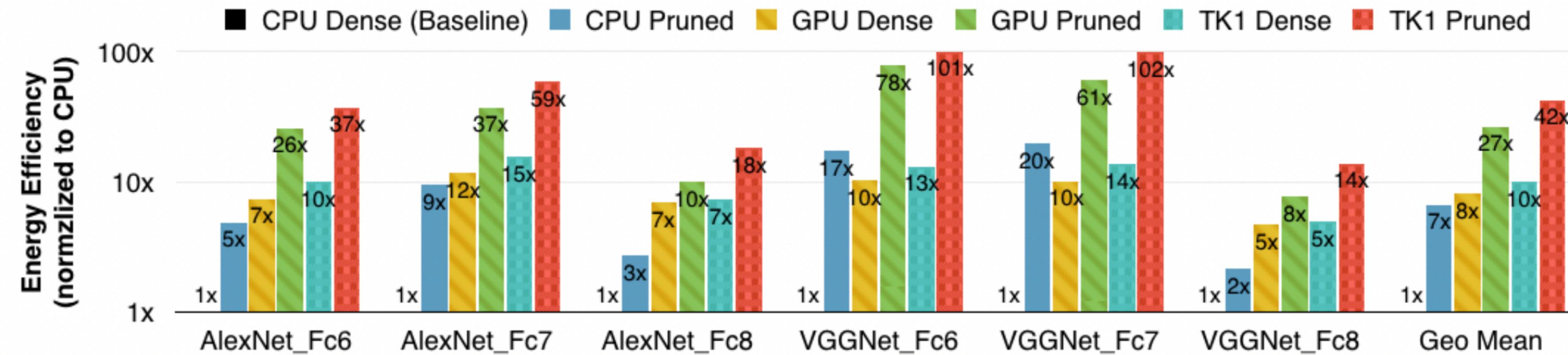
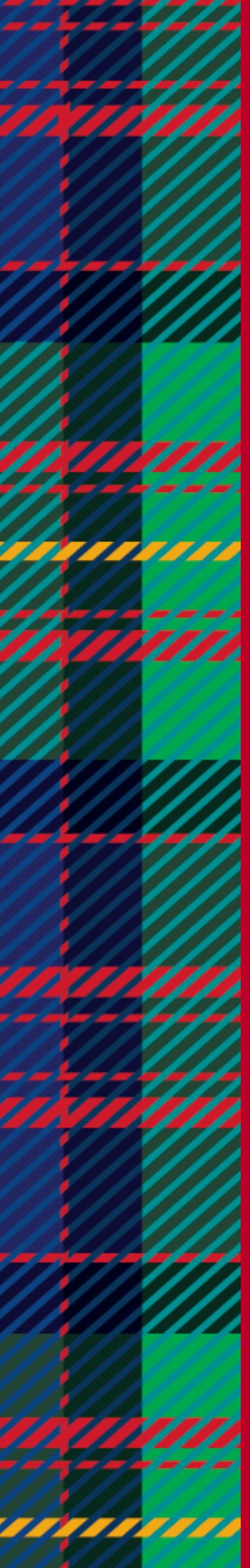


Figure 10: Compared with the original network, pruned network layer takes  $7\times$  less energy on CPU,  $3.3\times$  less on GPU and  $4.2\times$  less on mobile GPU on average. Batch size = 1 targeting real time processing. Energy number normalized to CPU.

**Thank you!**

# Questions

- In fact, the unstructured sparsity (e.g., pruning synapses) in this paper is a little bit hard to convert to real speedup. Is there some way to regularize the sparsity?
- In this paper, the authors set thresholds for each layer manually, which requires much human effort. Besides, the layer sparsity may also have some dependency. For example, you may prune the first layer 90% or the second layer 90% without loss of performance. However, you could not prune both of them 90%. Is there some way that could automatically prune the network as a whole to a targeted compression ratio while achieving the best accuracy?
- Is there some way that you could codesign your pruning and quantization with the system-level optimization (e.g., TVM)? For example, sometimes, you may find pruning more does not indicate lower latency due to some system optimization. Generally, multipliers of 8 are more friendly for system optimization.



# Carnegie Mellon University

## THE LOTTERY TICKET HYPOTHESIS: FINDING SPARSE, TRAINABLE NEURAL NETWORKS

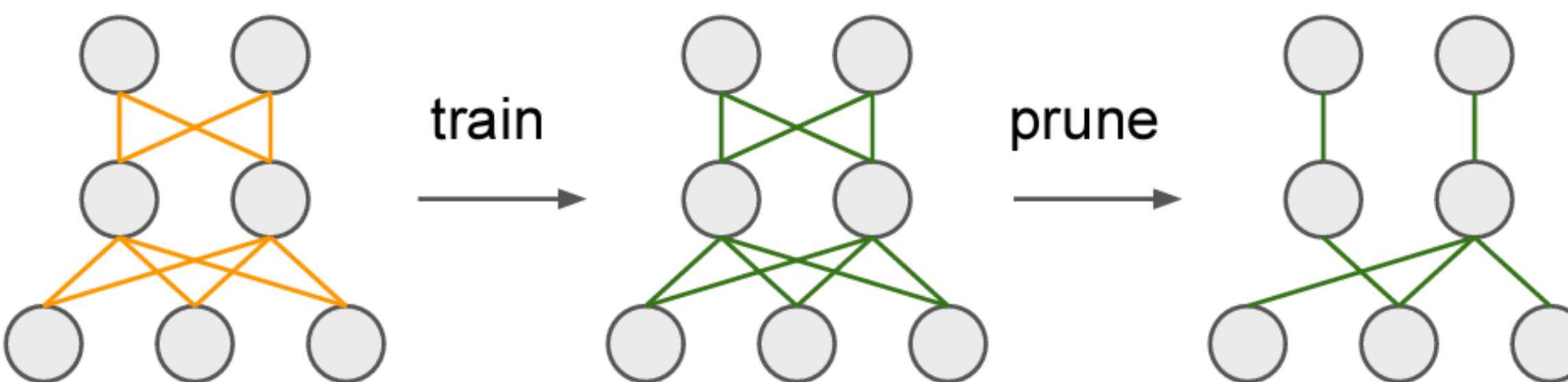
---

Authors: Jonathan Frankle (MIT), Michael Carbin (MIT)

Presenter: Rulin Shao, Machine Learning Department

# Motivating Question

We learned from the previous paper that model compression techniques like pruning can reduce the number of parameters by more than 90% without harming accuracy.



Then we may wonder:

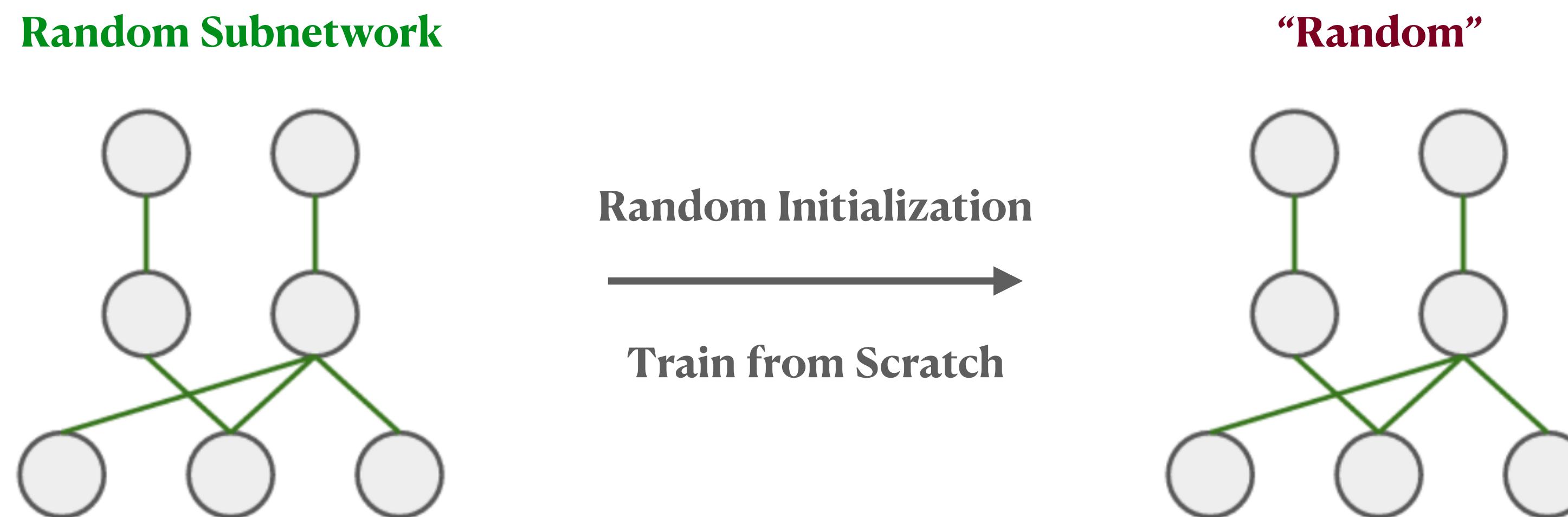
Since a smaller model can perform well on the task, can we train this sparse model from scratch to achieve comparable performance but with less computations?

LTH: Promising!

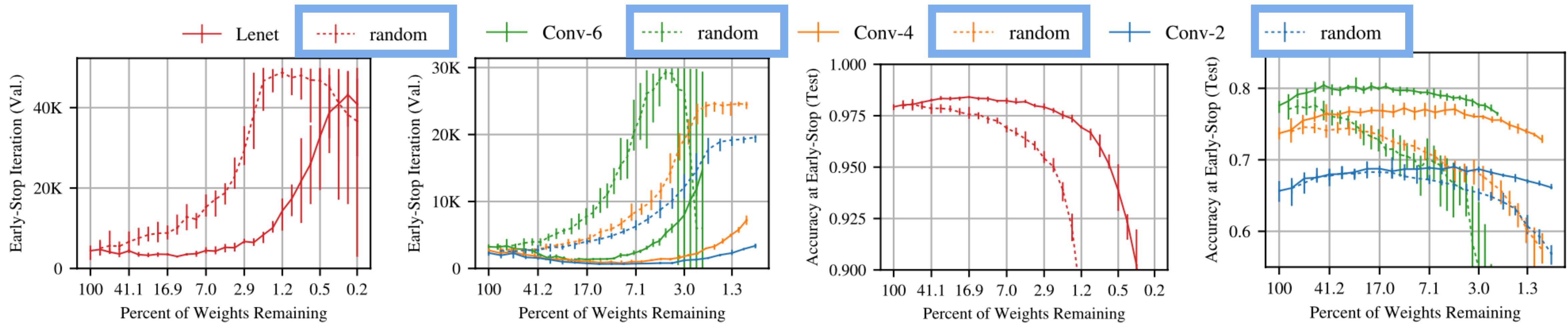
Carnegie  
Mellon  
University

# Proposal 1: Random Sparse Network?

What if we randomly sample a sparse subnetwork and train it from scratch?



# Proposal 1: Random Sparse Network?

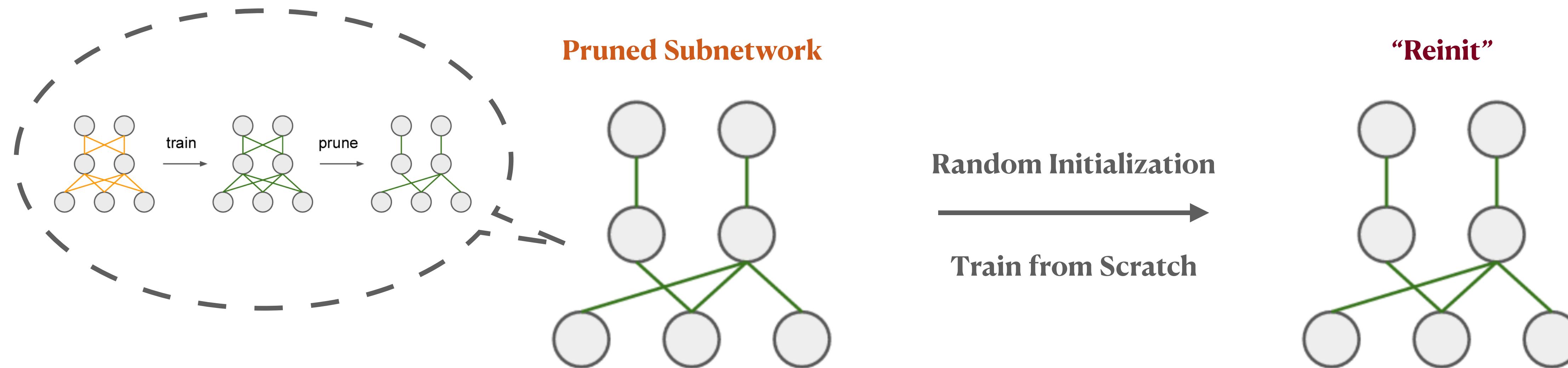


The **sparser** the network, the **slower** the learning and the **lower** the eventual test accuracy.

# Proposal 2: Topology Connections from Pruning?

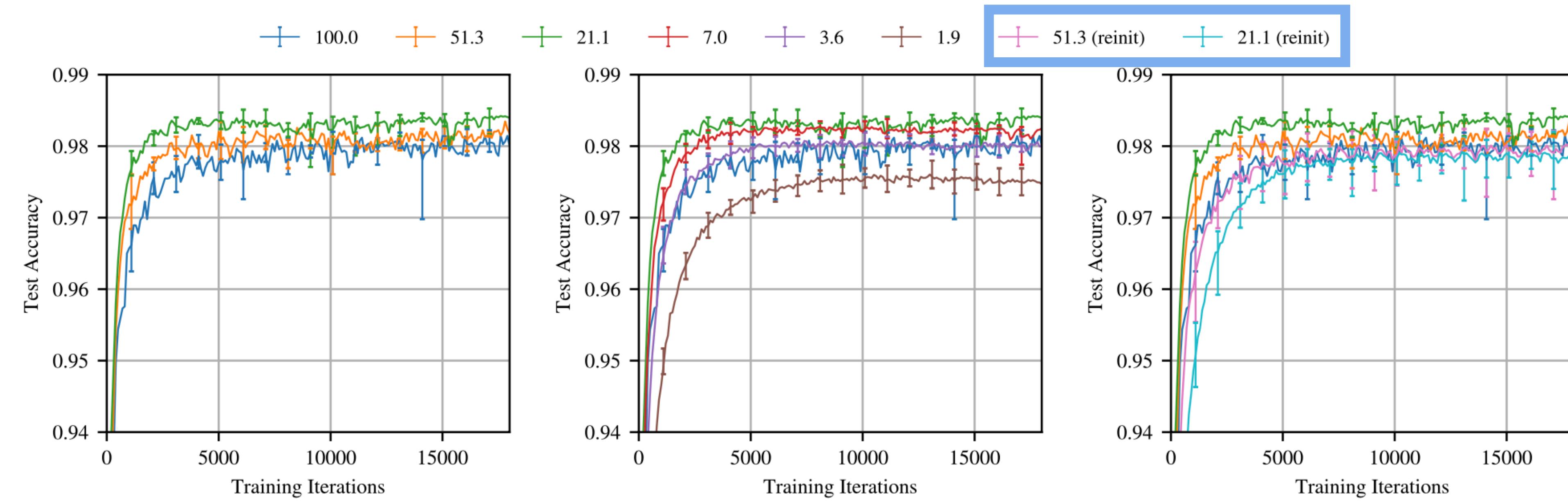
Well, if an arbitrary network doesn't work, what can we do next?

What if we take advantage of **the connection topology found by pruning**?



Carnegie  
Mellon  
University

# Proposal 2: Topology Connections from Pruning?



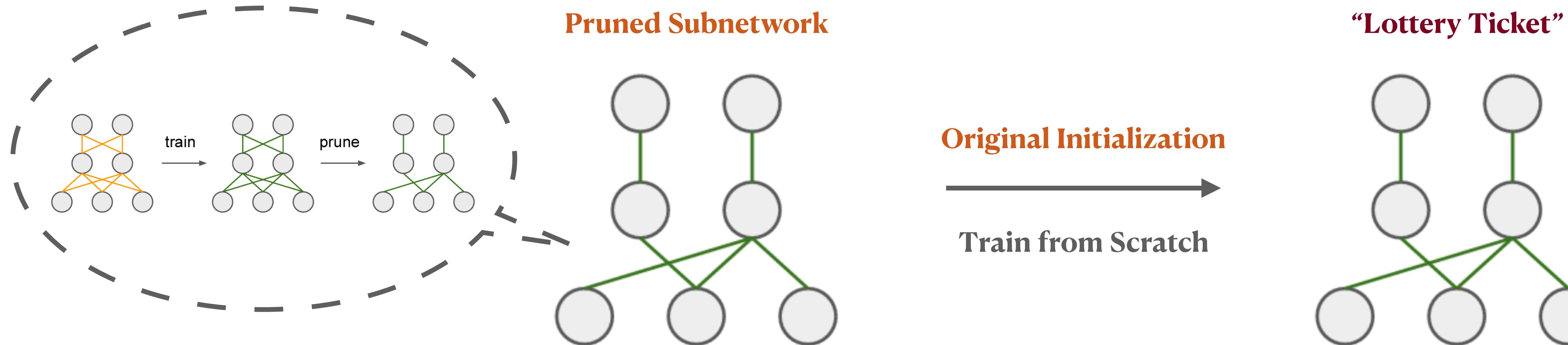
Not as good as the dense neural network.

**Structure alone cannot meet our goal.**

Carnegie  
Mellon  
University

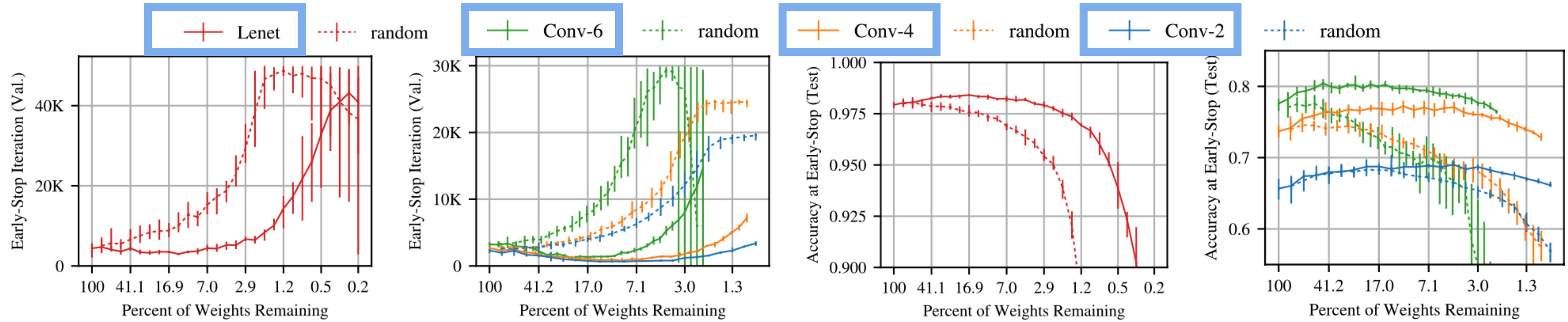
# Proposal 3: Topology Connections + Initialization?

What if we further preserve the original initialization?



Carnegie  
Mellon  
University

# Proposal 3: Topology Connections + Initialization?



Lottery tickets have **faster learning** and **higher test accuracy** than the dense network!

Carnegie  
Mellon  
University

We reached the goal!



Carnegie  
Mellon  
University

# Lottery Ticket Hypothesis (LTH)

Formally speaking:

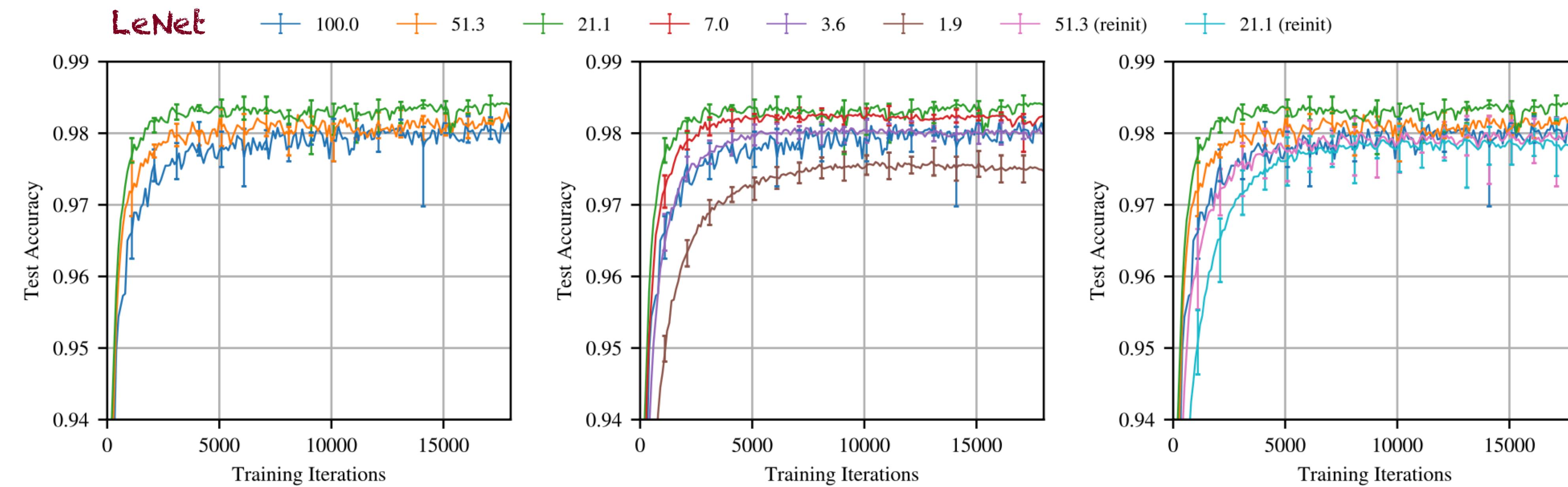
Found by unstructured pruning

A randomly-initialized, dense neural network contains a subnetwork that is initialized such that — when trained in isolation — it can match the test accuracy of the original network after training for at most the same number of iterations.

Using the original initialization

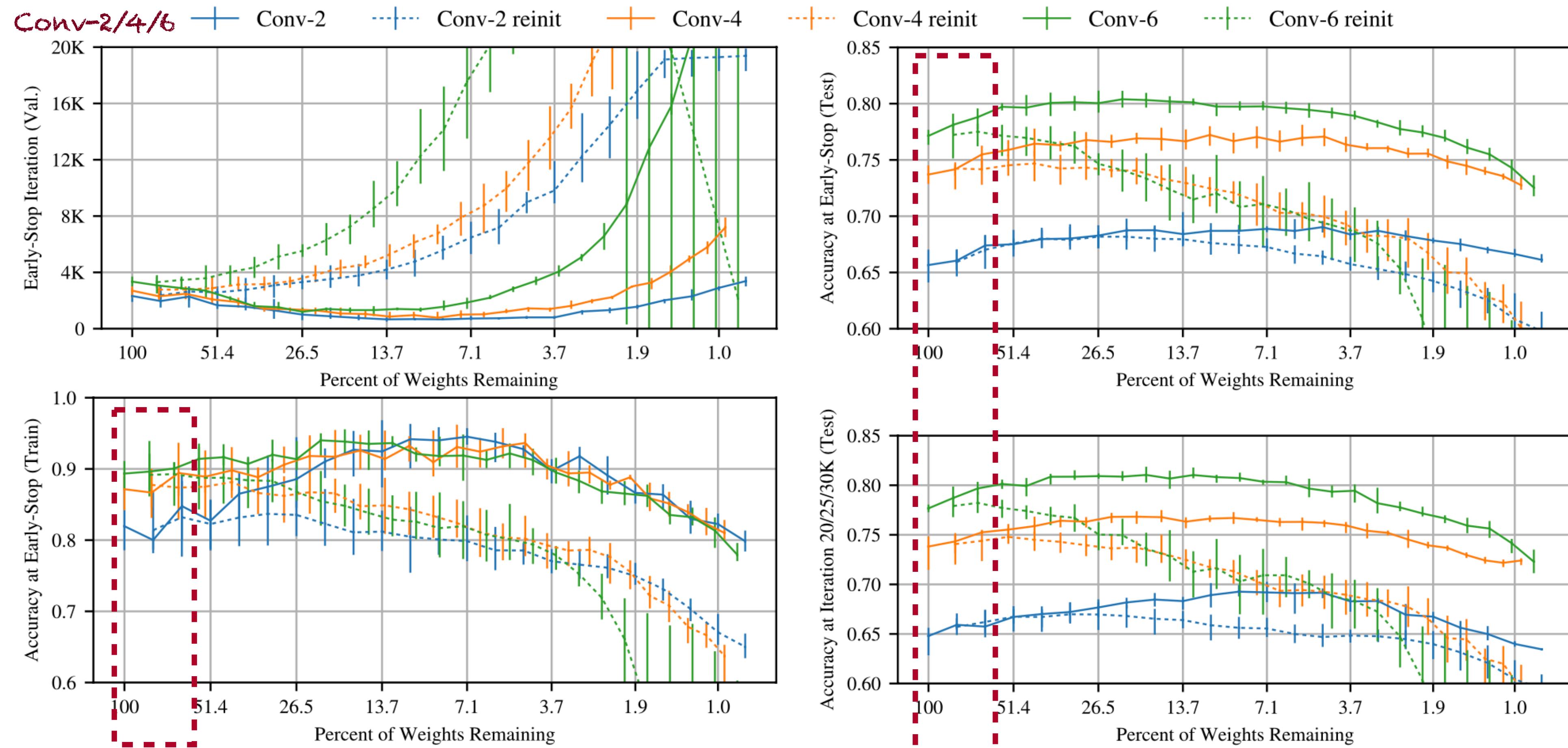
Carnegie  
Mellon  
University

# LTH: Experiments



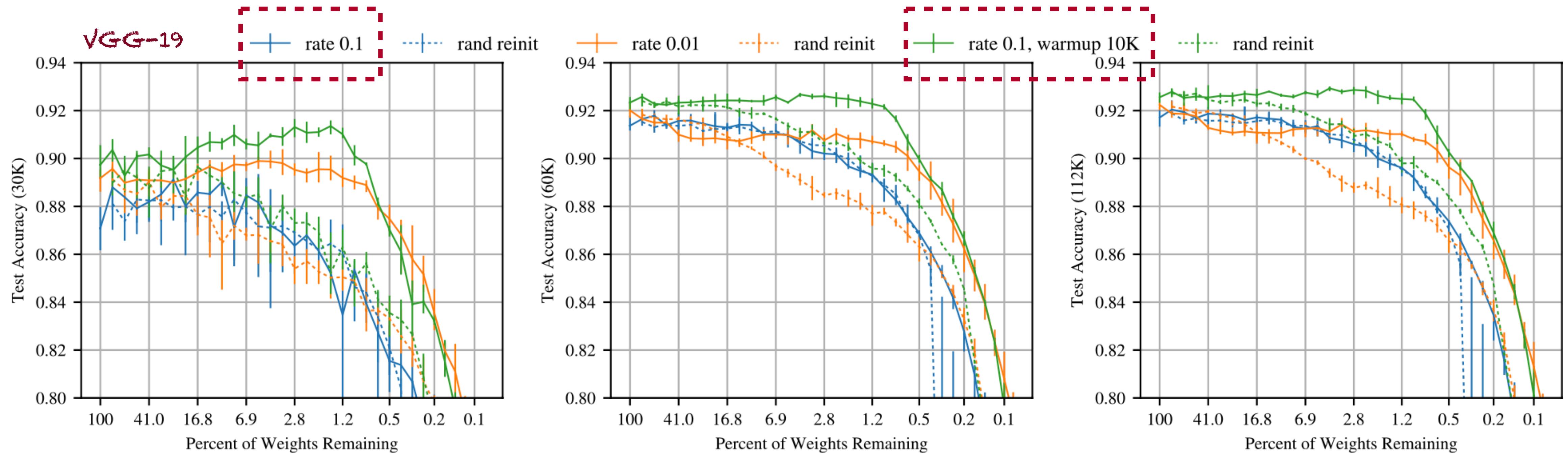
Carnegie  
Mellon  
University

# LTH: Experiments



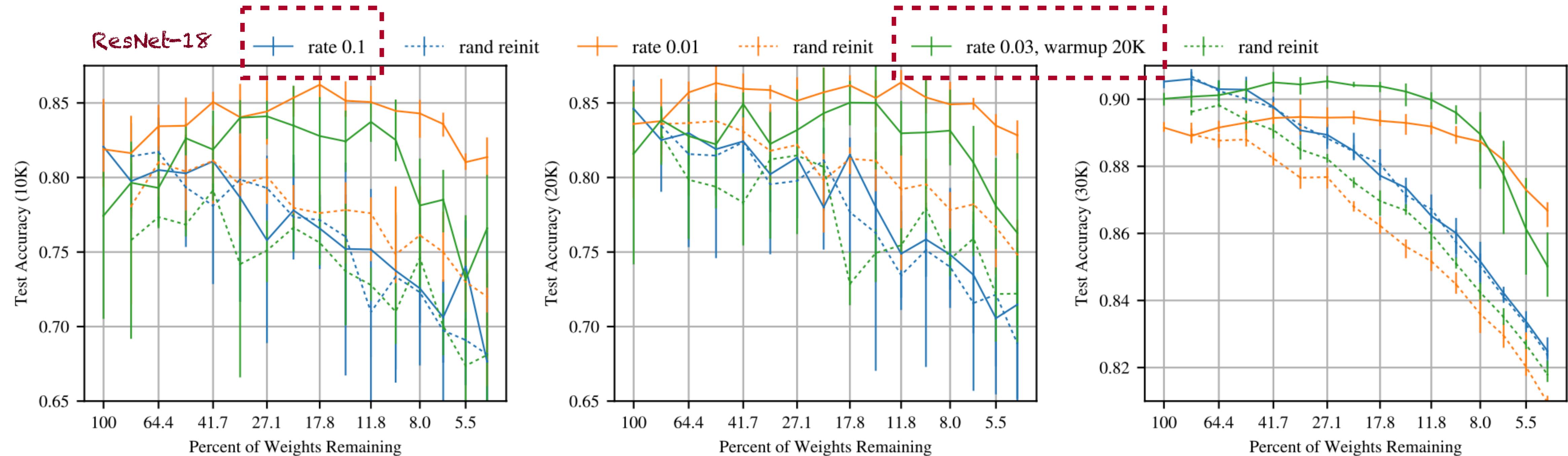
Carnegie  
Mellon  
University

# LTH: Experiments



Carnegie  
Mellon  
University

# LTH: Experiments



Carnegie  
Mellon  
University

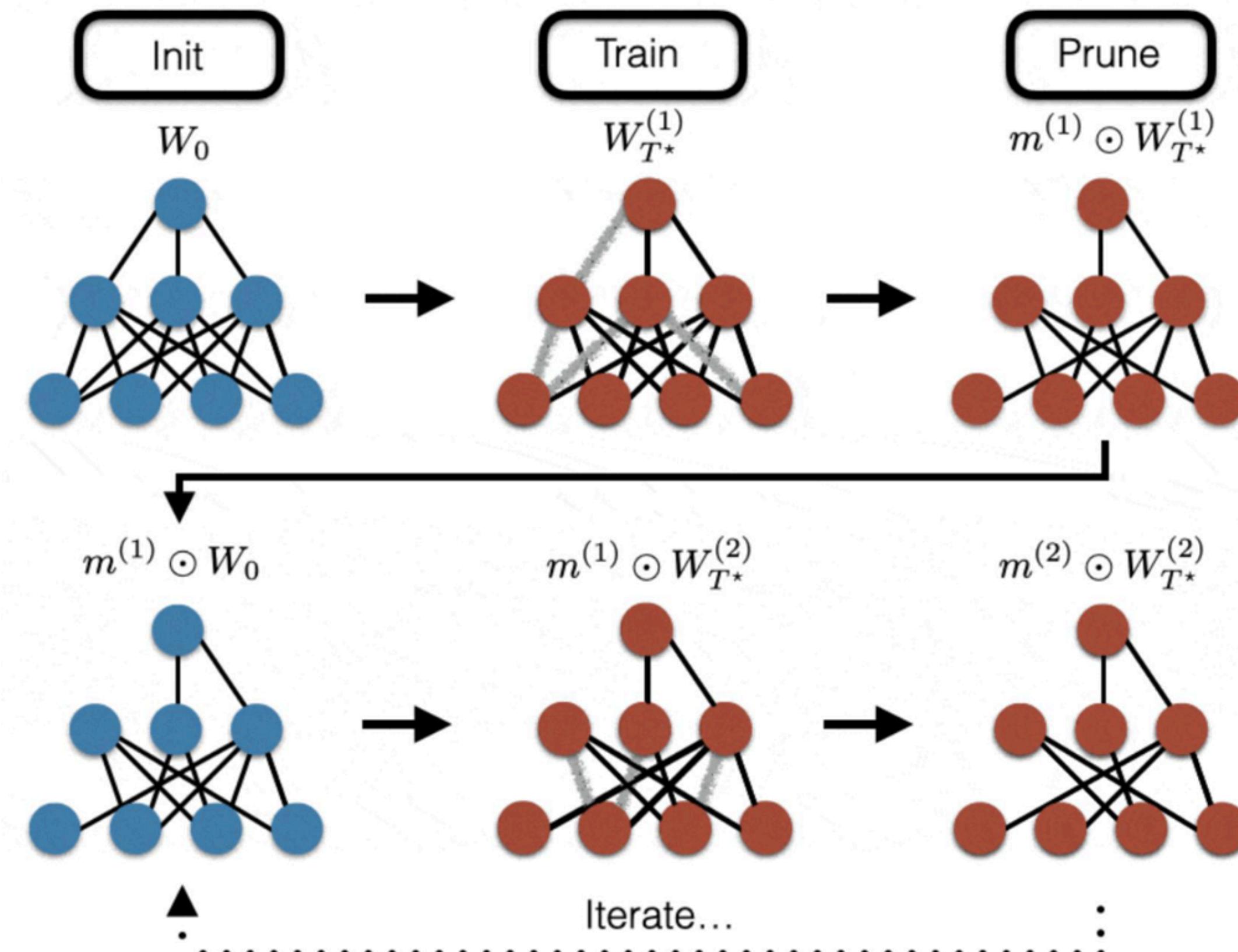
# LTH: one-shot v.s. iterative pruning

How does LTH find the winning tickets?

## — Pruning: one-shot v.s. iterative

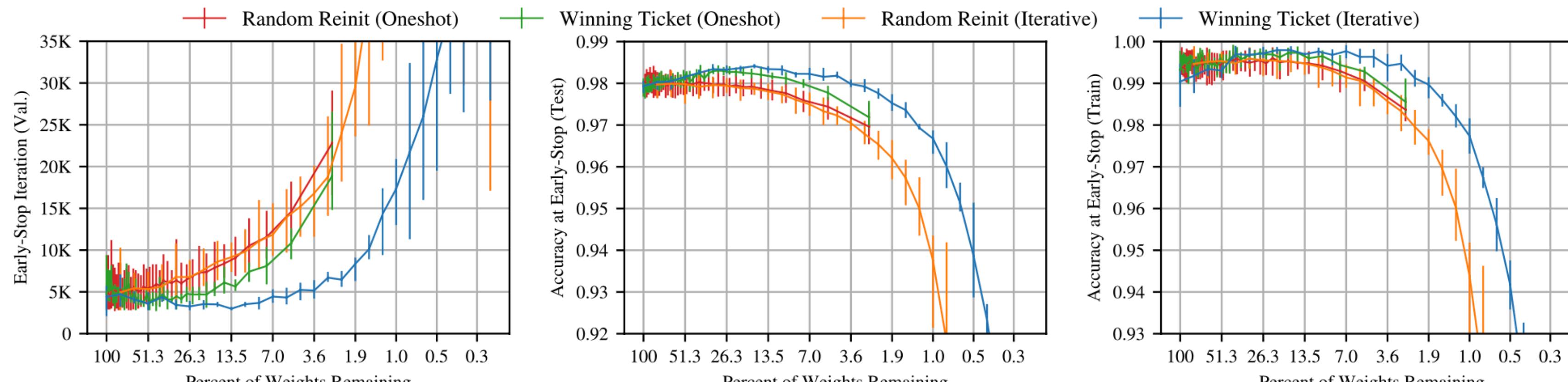
### Iterative Pruning:

1. Initialize the model
2. Train the model
3. Prune  $p^{\frac{1}{n}}\%$  of weights with smallest magnitudes
4. Reset values of remaining weights to their previous initial values
5. Repeat 2 - 4 until the desired sparsity is met

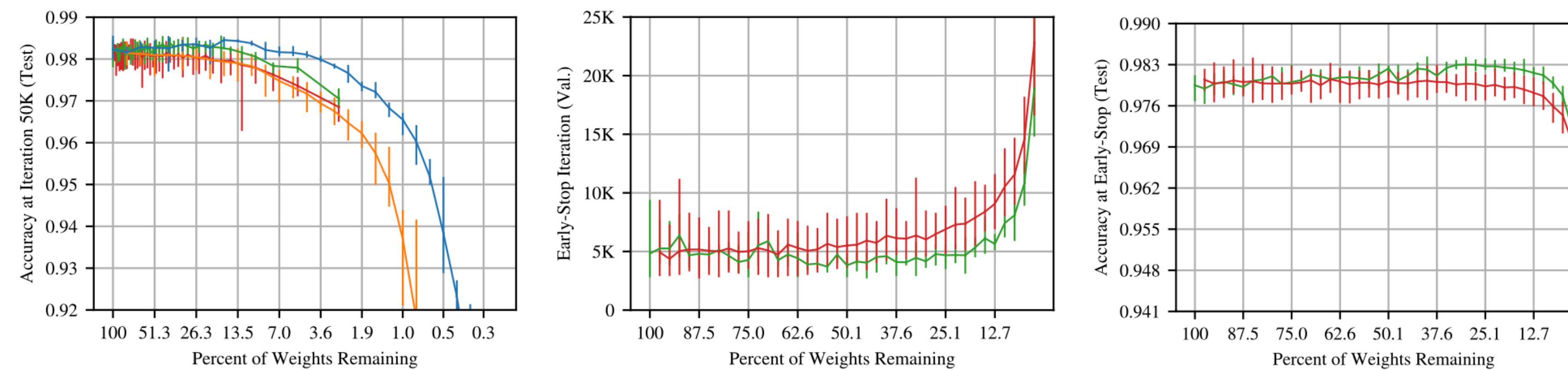


Carnegie  
Mellon  
University

# LTH: one-shot v.s. iterative pruning



(a) Early-stopping iteration and accuracy for all pruning methods.



(b) Accuracy at end of training.

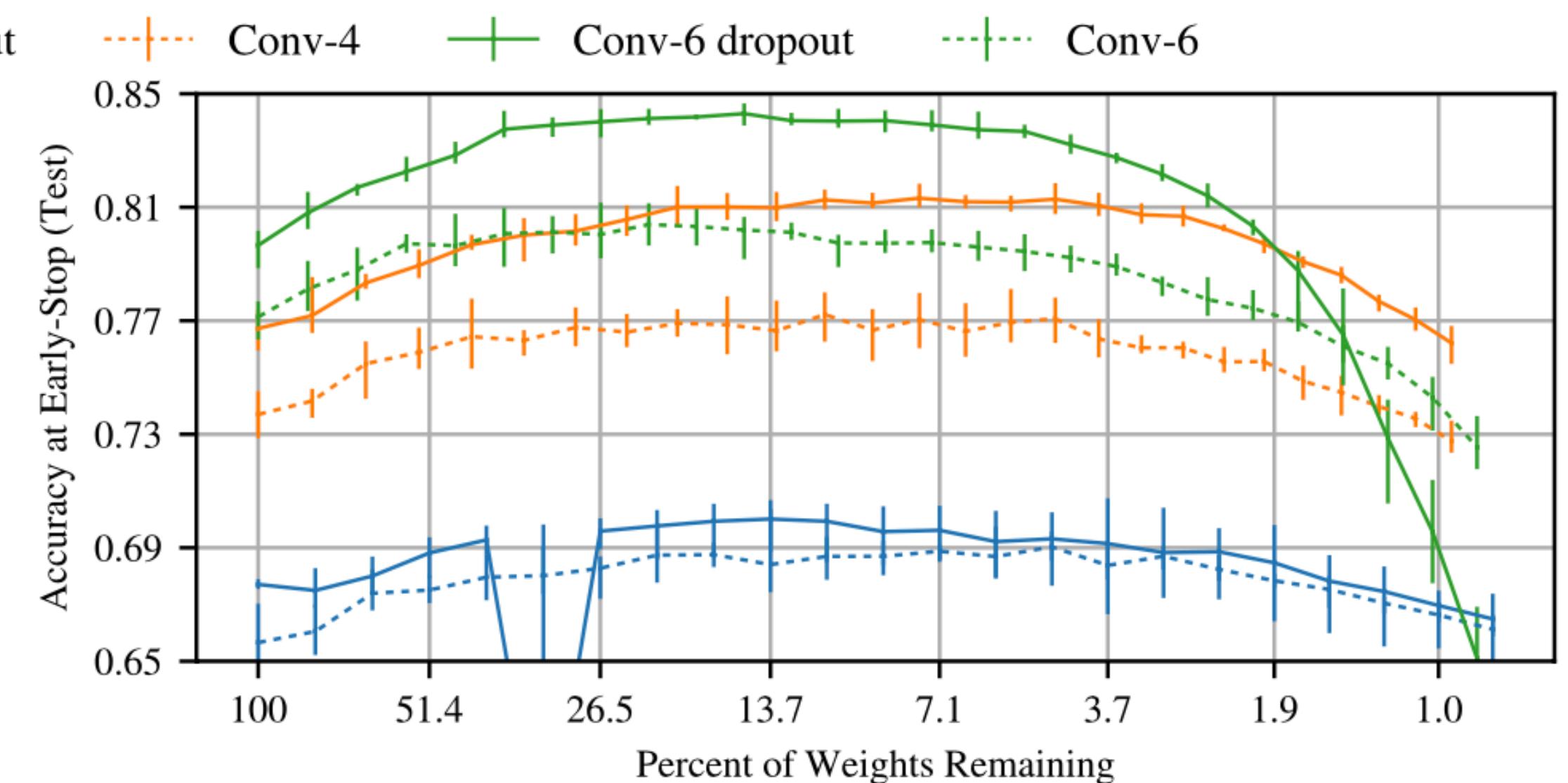
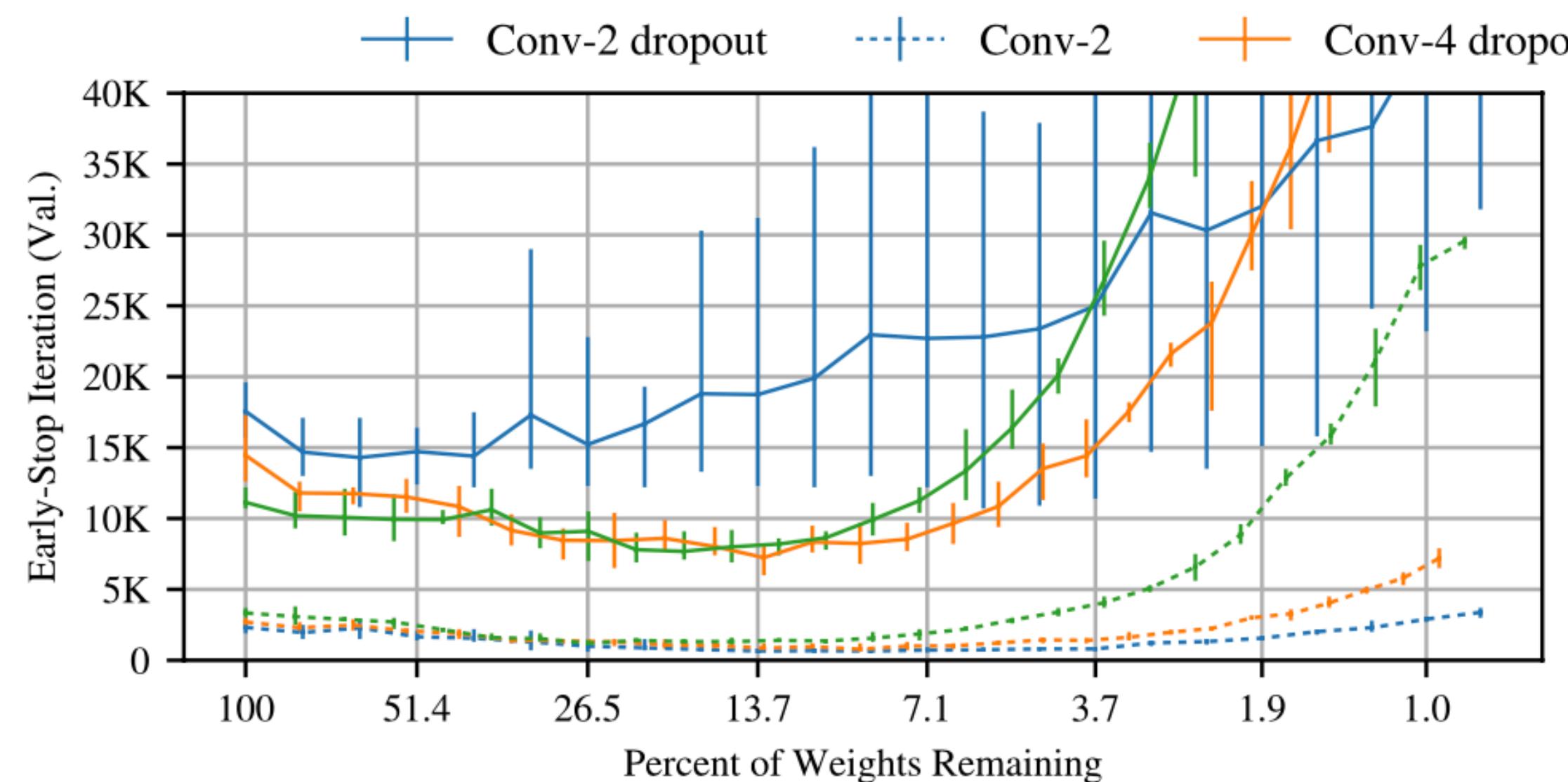
(c) Early-stopping iteration and accuracy for one-shot pruning.

# LTH: effect of dropout

Dropout is simultaneously training the ensemble of all subnetworks.

Dropout induce sparsity in the activations.

It is possible that dropout-induced sparsity primes a network to be pruned.



# LTH: Issues and Related Works

Have we really reached the goal? Not really...

- Finding the winning tickets (i.e., the pruning process) is too expensive. [2,3,4,5]
- Unstable: not useful for large learning rate and structured pruning (e.g.,  $l_1$ -norm pruning). [6]
- The original initialization contains the pre-training information. LT is not strictly training from scratch ...

Limitations:

- Only involved small detests (MNIST and CIFAR-10) [2,3,4,5]
- Only considered unstructured pruning. [6]
- The resulting architectures are not optimized for modern libraries or hardware. [5]

[1] Frankle, Jonathan, and Michael Carbin. "The lottery ticket hypothesis: Finding sparse, trainable neural networks." (ICLR 2019).

[2] Evci, Utku, et al. "Rigging the lottery: Making all tickets winners." (ICML 2020).

[3] Jayakumar, Siddhant, et al. "Top-kast: Top-k always sparse training." (NeurIPS 2020).

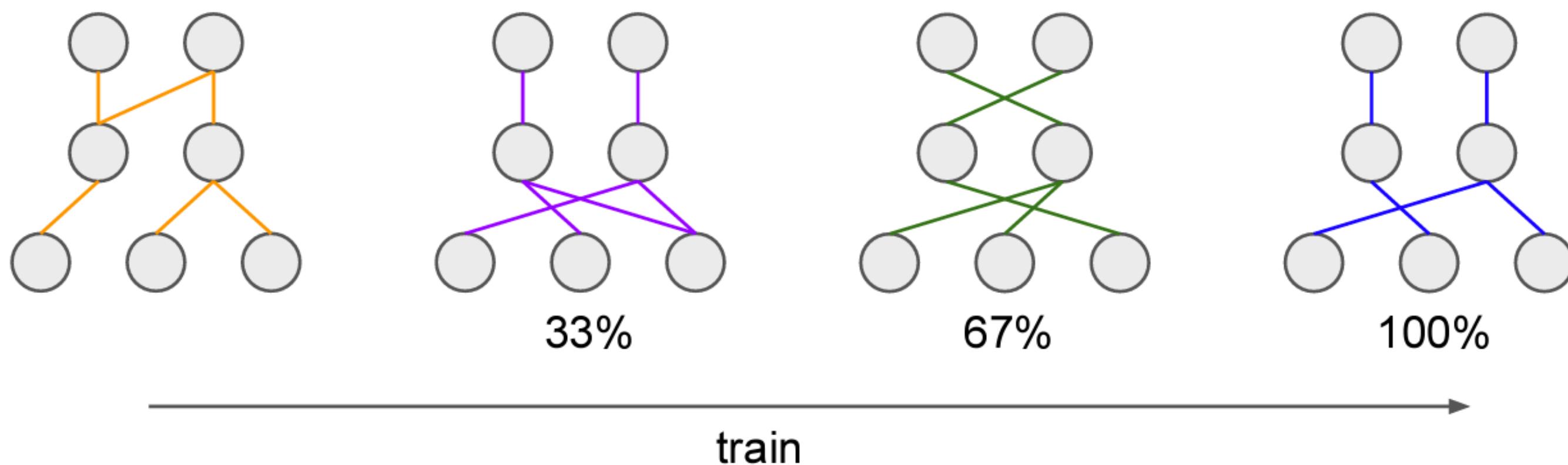
[4] Dettmers, Tim, and Luke Zettlemoyer. "Sparse networks from scratch: Faster training without losing performance." (2019).

[5] Zhou, Aojun, et al. "Learning n: m fine-grained structured sparse neural networks from scratch." (ICLR 2021).

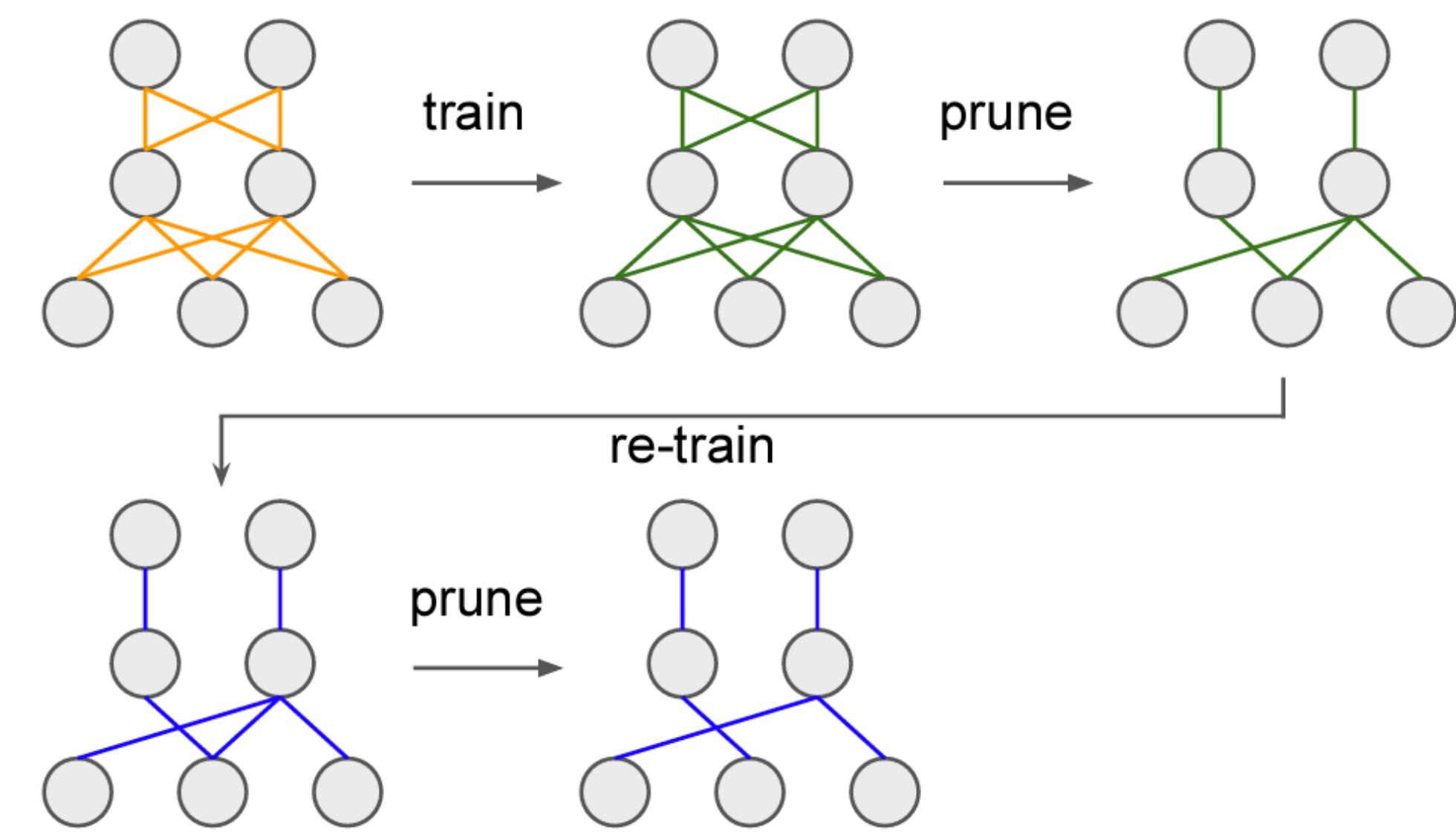
[6] Liu, Zhuang, et al. "Rethinking the value of network." (ICLR 2019)

# Follow-up Works: Sparse Training

## Sparse-to-sparse



## Dense-to-sparse

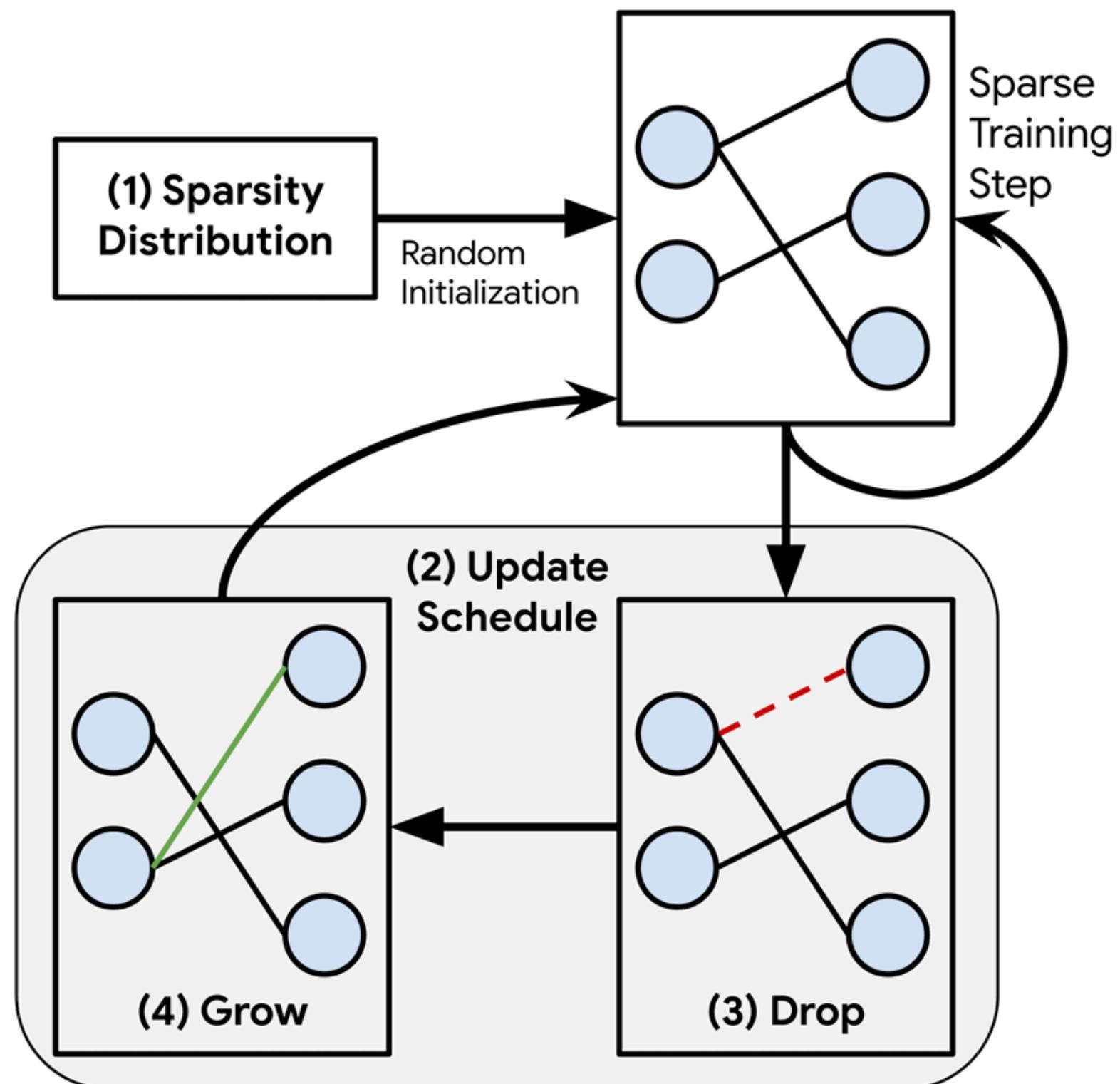


- Sparsity is maintained throughout the whole training process
- Connection topology can change during training
- Goal: Achieving the same performance as dense training with minimal computations

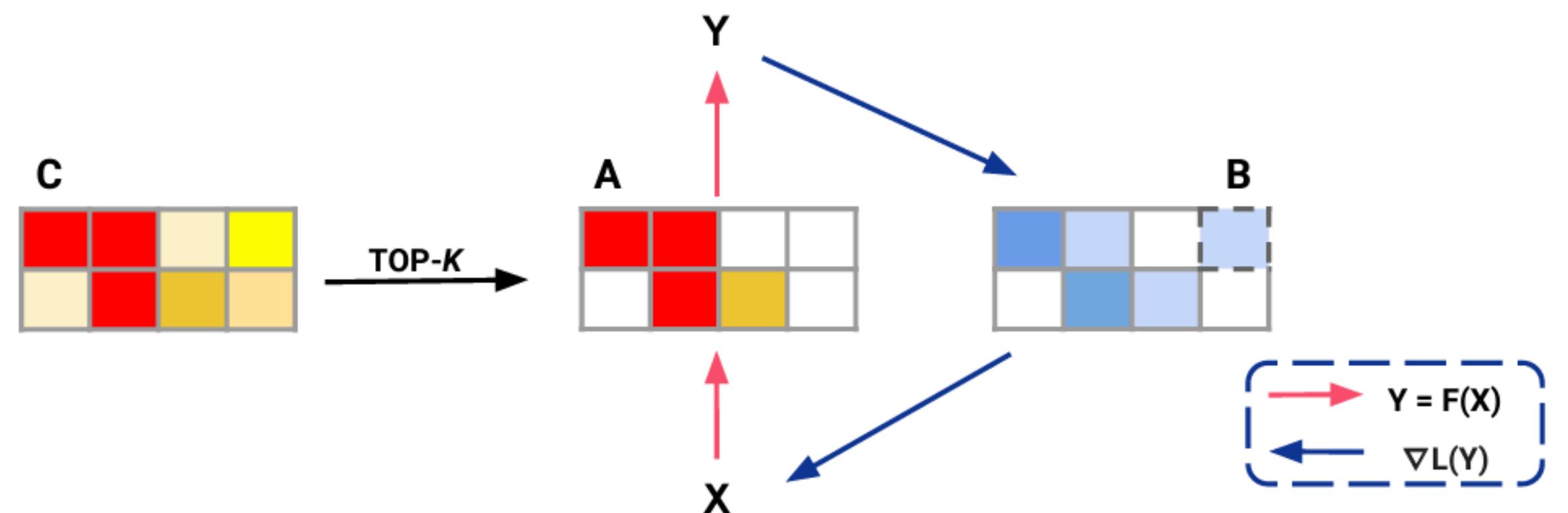
- (Iteratively) train the dense model, and then sparsify the “trained” dense model

# Follow-up Works: Sparse Training

## Rigging the lottery [1]



## Top-KAST [2]



**Sparse Momentum** [3]: redistribute removed weights using sparse momentum.  
**N:M** [4]: only N weights are non-zero for every continuous M weights.  
**Rethinking** [5]: consider structured pruning.

...

[1] Evci, Utku, et al. "Rigging the lottery: Making all tickets winners." (ICML 2020).

[2] Jayakumar, Siddhant, et al. "Top-kast: Top-k always sparse training." (NeurIPS 2020).

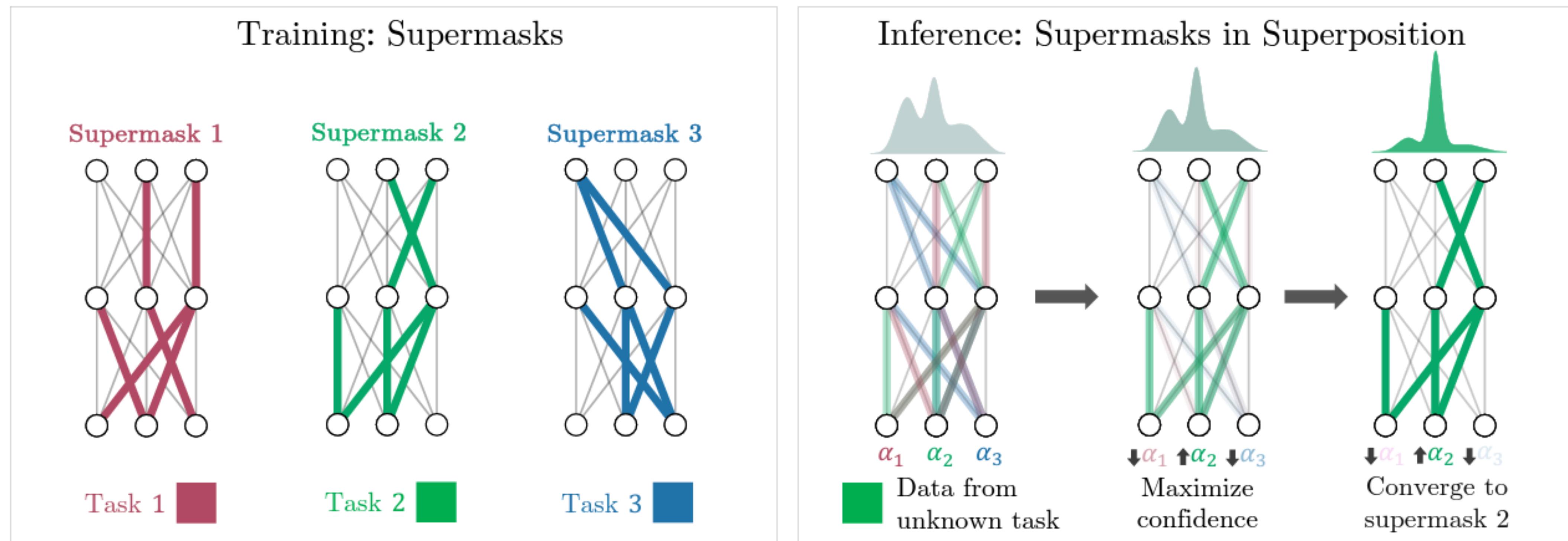
[3] Dettmers, Tim, and Luke Zettlemoyer. "Sparse networks from scratch: Faster training without losing performance." (2019).

[4] Zhou, Aojun, et al. "Learning n: m fine-grained structured sparse neural networks from scratch." (ICLR 2021).

[5] Liu, Zhuang, et al. "Rethinking the value of network." (ICLR 2019)

# Structural Information in Subnetworks

**SupSup [1]:** The base network is randomly initialized, learn supermasks for different tasks.



# Lottery Ticket Hypothesis Conjecture

(Untested) **Conjecture:** SGD **seeks out and trains a subset of well-initialized weights.** Dense, randomly-initialized networks are easier to train than the sparse networks that result from pruning because **there are more possible subnetworks from which training might recover a winning ticket.**

# One (Possible) Perspective to Support the Conjecture: Equivalences Between Sparse Models and Neural Networks [1]

The training process of neural networks (e.g. with weight decay or early stopping) are “hunting” for a sparse representation in terms of its learned features!

- Equivalences can be derived between sparse models and neural networks.
- Weight decay and early stopping can both be sparse-inducing in neural networks.

# Another Perspective: Equivalences Between Sparse Models and Neural Networks

The simplest case: The lasso is equivalent to a two-layer neural network\* fit with weight decay:

Lasso		A two-layer neural network fit with weight decay
$\underset{\beta}{\text{minimize}} \quad \sum_{i=1}^n (y_i - x_i^\top \beta)^2 + 2\lambda \sum_{j=1}^d  \beta_j $	$\longleftrightarrow$	$\underset{u,v}{\text{minimize}} \quad \sum_{i=1}^n \left( y_i - \sum_{j=1}^d x_{ij} u_j v_j \right)^2 + \lambda \sum_{j=1}^d (u_j^2 + v_j^2)$ <p style="color: #800000; font-size: 0.8em; margin-top: 10px;">Optimize in the 2d dimension (instead of d)</p>

$$\forall c, \underset{ab=c}{\text{min}}(a^2 + b^2) = 2|c| \text{ is achieved at } a = b = \sqrt{|c|}$$

Carnegie  
Mellon  
University

\* A two-layer neural network with linear activation functions, no bias terms, and a very simple connectivity structure

# Another Perspective: Equivalences Between Sparse Models and Neural Networks

A  $k$ -layer neural network\* that has otherwise the same structure is in turn equivalent to an  $l_p$ -penalized regression problem, where  $p = 2/k < 1$ .

**An  $l_p$ -penalized regression problem**

$$\underset{\beta}{\text{minimize}} \quad \sum_{i=1}^n L\left(y_i, \sum_{j=1}^d x_{ij}\beta_j\right) + k\lambda \sum_{j=1}^d |\beta_j|^{2/k}$$



**K-layer neural network**

$$\underset{w^{(\ell)}, \ell=1,\dots,k}{\text{minimize}} \quad \sum_{i=1}^n L\left(y_i, \sum_{j=1}^d x_{ij} \prod_{\ell=1}^k w_j^{(\ell)}\right) + \lambda \sum_{j=1}^d \sum_{\ell=1}^k (w_j^{(\ell)})^2$$

\* A  $k$ -layer neural network with linear activation functions, no bias terms, and a very simple connectivity structure since the third layer

# Another Perspective: Equivalences Between Sparse Models and Neural Networks

Similarly observations can be extended to k-layer group-sparse models\* and (2 or 3-layer) fully connected NN with linear/ReLU activations... (Equivalent to a **group lasso problem**)

**A 2-layer fully connected NN with ReLU activation**

$$\underset{b^{(1)}, w^{(1)}, w^{(2)}}{\text{minimize}} \quad \sum_{i=1}^n L\left(y_i, \sum_{j=1}^{d_2} \phi(x_i^\top w_j^{(1)} + b_j^{(1)}) w_j^{(2)}\right) + \lambda \left(\sum_{j=1}^{d_2} \|w_j^{(1)}\|_2^2 + \|w^{(2)}\|_2^2\right)$$

**A group lasso problem**

$$\underset{\alpha, \beta, s^{(2)} \in \{-1, 1\}^{d_2}}{\text{minimize}} \quad \sum_{i=1}^n L\left(y_i, \sum_{j=1}^{d_2} \phi(x_i^\top \beta_j + \gamma_j) s_j^{(2)}\right) + 2\lambda \sum_{j=1}^{d_2} \|\beta_j\|_2.$$

\* A k-layer neural network with linear activation functions, no bias terms, and a very simple connectivity structure since the third layer

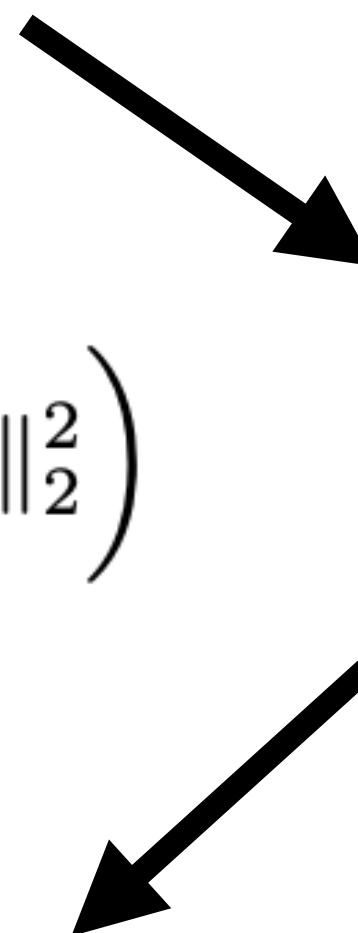
# Another Perspective: Equivalences Between Sparse Models and Neural Networks

**A group lasso problem**

$$\underset{\alpha, \beta, s^{(2)} \in \{-1, 1\}^{d_2}}{\text{minimize}} \sum_{i=1}^n L\left(y_i, \sum_{j=1}^{d_2} \phi(x_i^\top \beta_j + \gamma_j) s_j^{(2)}\right) + 2\lambda \sum_{j=1}^{d_2} \|\beta_j\|_2.$$

**A 2-layer fully connected NN with ReLU activation**

$$\underset{b^{(1)}, w^{(1)}, w^{(2)}}{\text{minimize}} \sum_{i=1}^n L\left(y_i, \sum_{j=1}^{d_2} \phi(x_i^\top w_j^{(1)} + b_j^{(1)}) w_j^{(2)}\right) + \lambda \left(\sum_{j=1}^{d_2} \|w_j^{(1)}\|_2^2 + \|w^{(2)}\|_2^2\right)$$



Due to what we know about the group lasso penalty, many of the parameters  $\beta_j$  will be sparse at the solution

The same will be true for the feature map:

$$x \mapsto \phi(x_i^\top w_j^{(1)} + b_j^{(1)}) w_j^{(2)}, \quad j = 1, \dots, d_2.$$

That is, many of these functions will be either the zero map, or a constant map .

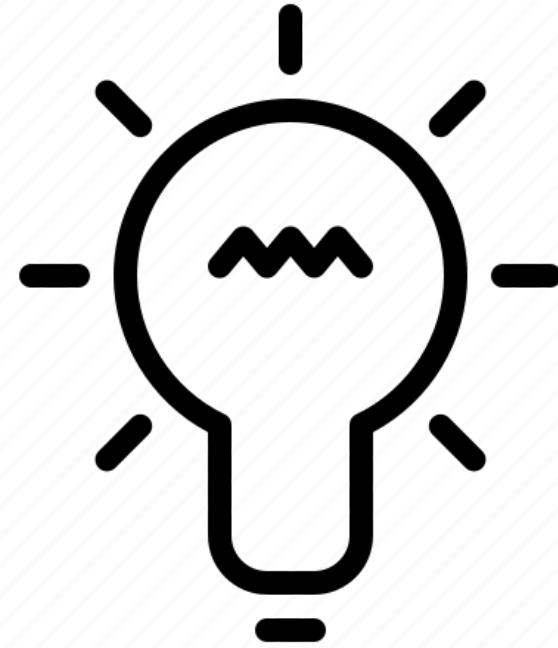
# Another Perspective: Equivalences Between Sparse Models and Neural Networks

$$x \mapsto \phi(x_i^\top w_j^{(1)} + b_j^{(1)}) w_j^{(2)}, \quad j = 1, \dots, d_2.$$

That is, many of these functions will be either the zero map, or a constant map .

**Therefore, a standard two-layer neural network, trained with weight decay, is actually “hunting” for a sparse representation in terms of its learned features!**

# Another Perspective: Equivalences Between Sparse Models and Neural Networks



One possible takeaway from this perspective:

If the neural networks are learning sparse representation themselves, it shouldn't be surprising that we can prune a learned network to a sparse structure.

But in lower dimensions, we need to rely on **good initializations/appropriate learning rate/careful construction of the model structure** ... for current optimizers to find a good local optimal.

E.g. rewinding[2,3]...

- These subnetworks only reach full accuracy when they are stable to SGD noise, which occurs at initialization in some settings but only after some amount of training in others. [2]
- Weight rewinding, learning rate rewinding[3] ...

[1] Tibshirani, Ryan J. "Equivalences Between Sparse Models and Neural Networks." (2021).

[3] Lezama, José, et al. "Ole: Orthogonal low-rank embedding-a plug and play geometric loss for deep learning." (CVPR 2018)

[4] Wright, John, and Yi Ma. High-dimensional data analysis with low-dimensional models: Principles, computation, and applications. Cambridge University Press, 2022.

# Discussion Questions

1. Why is it hard for SGD to learn sparse representations from scratch? Why initialization is important to SGD (e.g. thinking of rewinding[1,2] techniques...)?
2. How to make the lottery tickets fit better for modern libraries (possibly at the cost of being less sparse)?
3. Can we use graphical models to learn the sparse connection topology?
4. Can we reformulate the sparse training problem as learning low-rank representations for high-dimensional data? E.g. using nuclear norm (OLE loss[3])/rate distortion (MCR loss[4]) to induce more interpretable sparse models instead of heuristically searching for sparse connections using pruning?

[1] Frankle, Jonathan, et al. "Linear mode connectivity and the lottery ticket hypothesis." (*ICML 2020*)

[2] Renda, Alex, Jonathan Frankle, and Michael Carbin. "Comparing rewinding and fine-tuning in neural network pruning." (*ICLR 2020*)

[3] Lezama, José, et al. "Ole: Orthogonal low-rank embedding-a plug and play geometric loss for deep learning." (*CVPR 2018*)

[4] Wright, John, and Yi Ma. High-dimensional data analysis with low-dimensional models: Principles, computation, and applications. Cambridge University Press, 2022.

# Discussion Questions

1. Why is it hard for SGD to learn sparse representations from scratch?
  - These subnetworks only reach full accuracy when they are stable to SGD noise, which occurs at initialization in some settings but only after some amount of training in others. [1]
  - Weight rewinding, learning rate rewinding[2] ...
2. How to make the lottery tickets fit better for modern libraries (possibly at the cost of being less sparse)?
  - M:N, Low-rank decomposition constrain ...
3. Can we use graphical models to learn the sparse connection topology?
4. Can we reformulate the sparse training problem as learning low-rank representations for high-dimensional data? E.g. using nuclear norm (OLE loss[3])/rate distortion (MCR loss[4]) to induce more interpretable sparse models instead of heuristically searching for sparse connections using pruning?

[1] Frankle, Jonathan, et al. "Linear mode connectivity and the lottery ticket hypothesis." (*ICML 2020*)

[2] Renda, Alex, Jonathan Frankle, and Michael Carbin. "Comparing rewinding and fine-tuning in neural network pruning." (*ICLR 2020*)

[3] Lezama, José, et al. "Ole: Orthogonal low-rank embedding-a plug and play geometric loss for deep learning." (*CVPR 2018*)

[4] Wright, John, and Yi Ma. High-dimensional data analysis with low-dimensional models: Principles, computation, and applications. Cambridge University Press, 2022.