

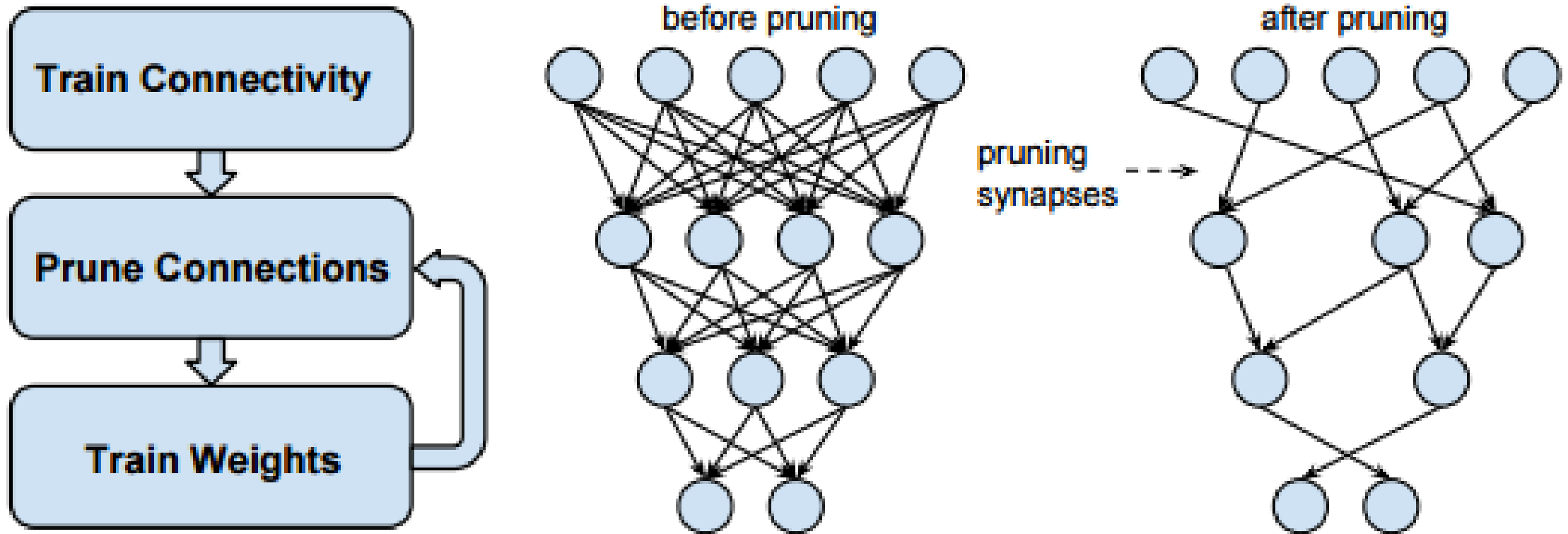
Exploring synapse pruning for energy-efficiency in Neural Networks (NN)

Aayush Ankit, Chankyu Lee, Priyadarshini Panda
ECE, Purdue University

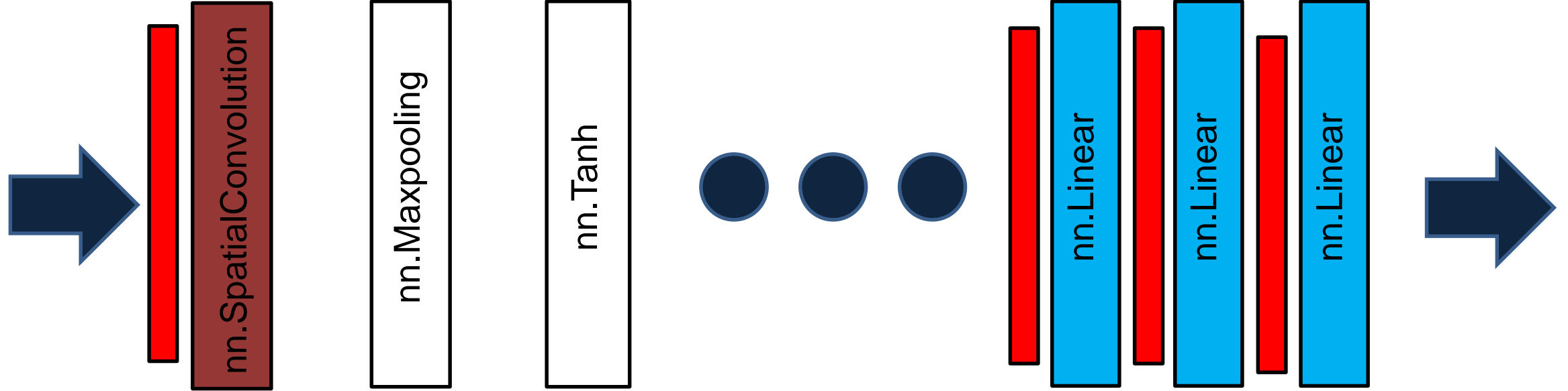
Motivation

- General purpose computing systems involve frequent and high bandwidth data transfers between memory and the computation core.
- Typical data flow :- DRAM – SRAM – Computation Core.
- **DRAM fetches extremely energy hungry.** Hence need efficient ways to store data in the memory.
- **Pruning exposes the sparsity in the NN** thereby, providing **avenues for data compression.**
- Data Compression can be significant driven (low overhead) for example – Frequent Pattern/Value compression, other encoding schemes etc.

What is pruning a network ?



Our Implementation



`nn.Prune`



Weight

0.9	0.8	0.1
0.2	0.3	0.7
0.6	0.4	0.5



Mask

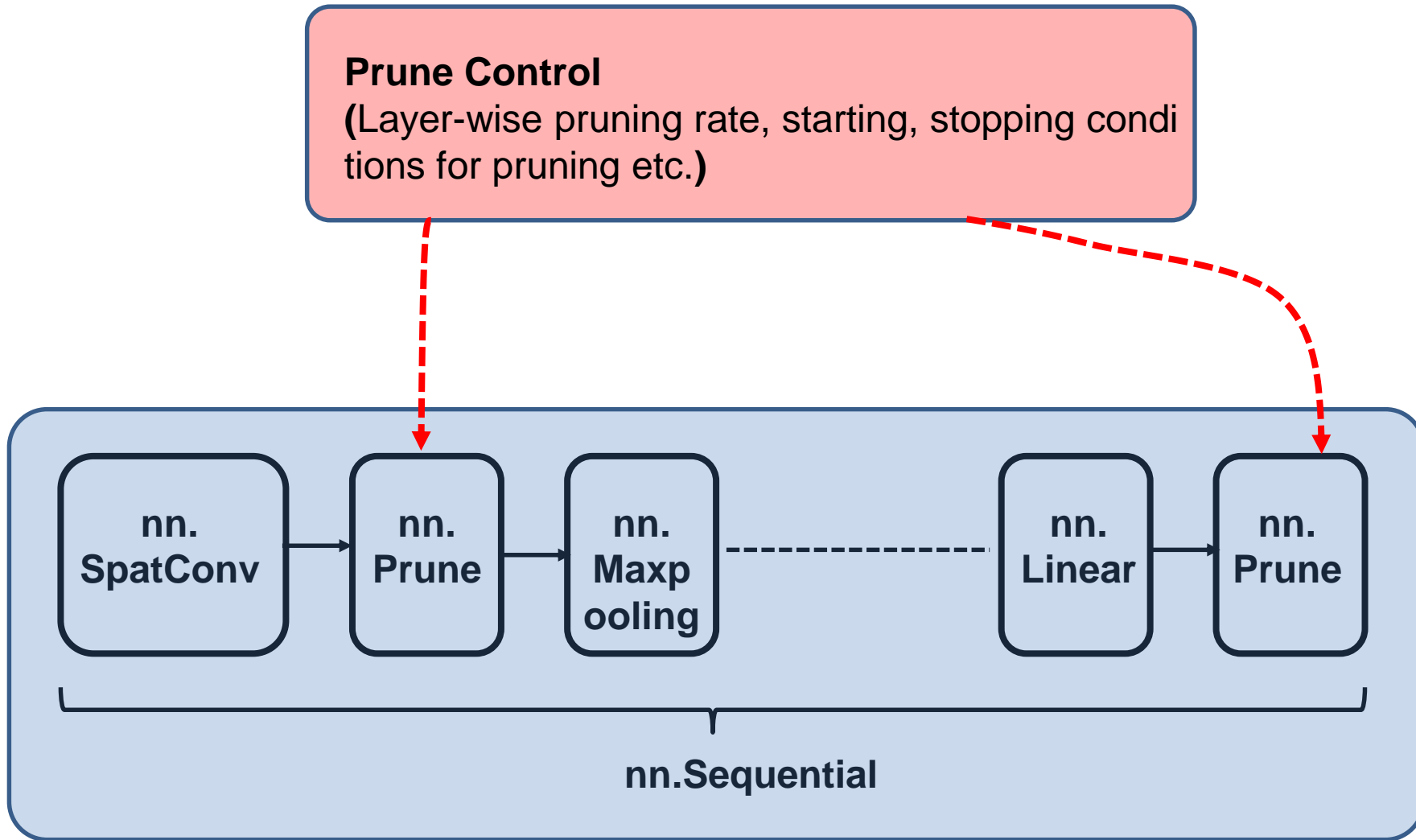
1	1	0
0	0	1
1	0	1



Pruned Weight

0.9	0.8	0
0	0	0.7
0.6	0	0.5

Our Implementation (Continued)



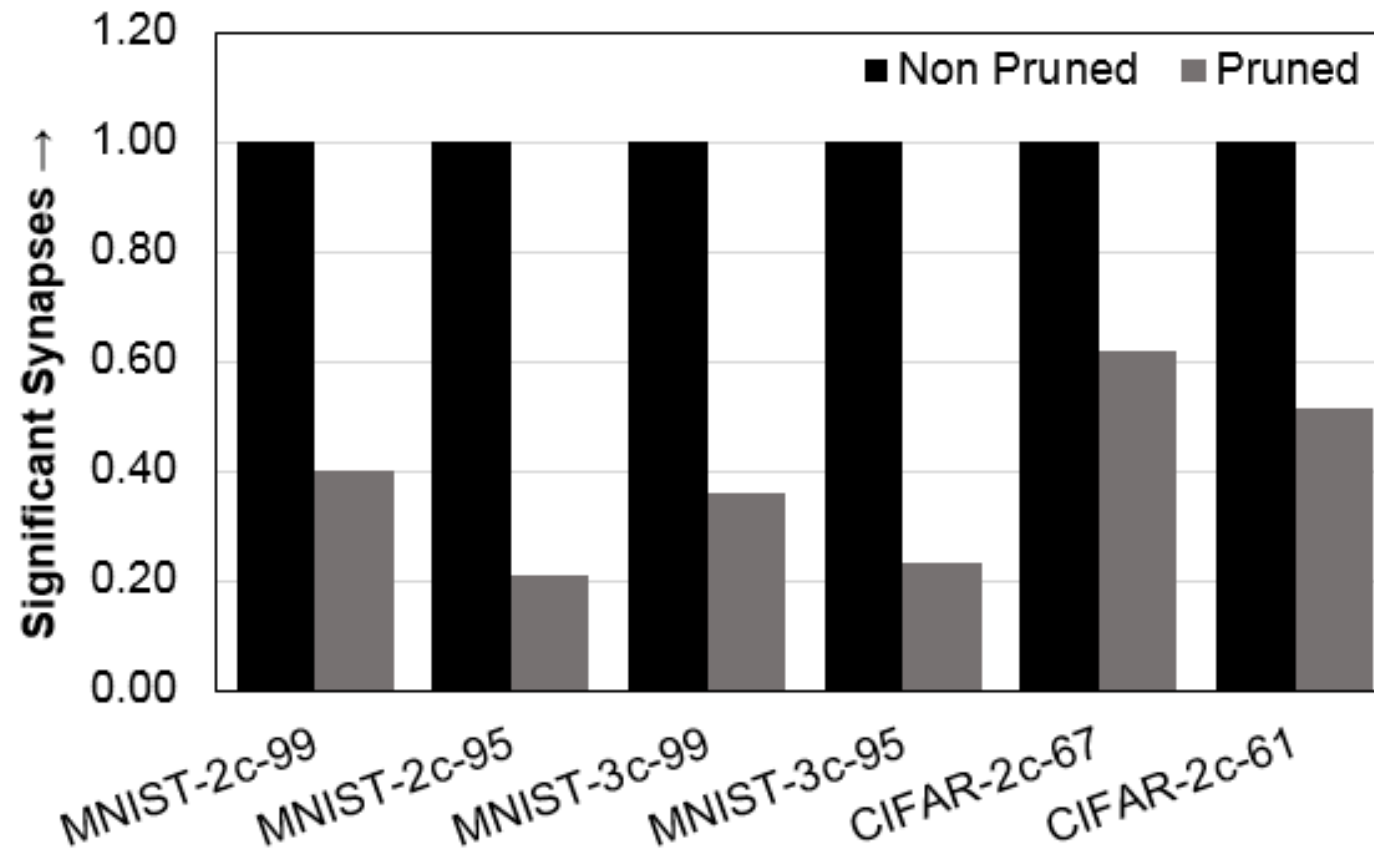
Experimental Methodology

➤ The “**nn.Prune**” functionality was implemented using masks to modify the “**forward**” and “**backward**” function in **Torch**.

➤ Simulation benchmarks :-

Dataset	Neural Network (NN) structure
MNIST	28x28 – 6c5 – 16c5 – 120 – 84 – 10
MNIST	28x28 – 6c3 – 16c3 – 36c4 – 240-120 – 84 – 10
CIFAR-10	3x32x32 – 6c5 – 16c5 – 120 – 84 – 10
CIFAR-10	3x32x32 – 6c3 – 16c3 – 36c4 – 240-120 – 84 – 10

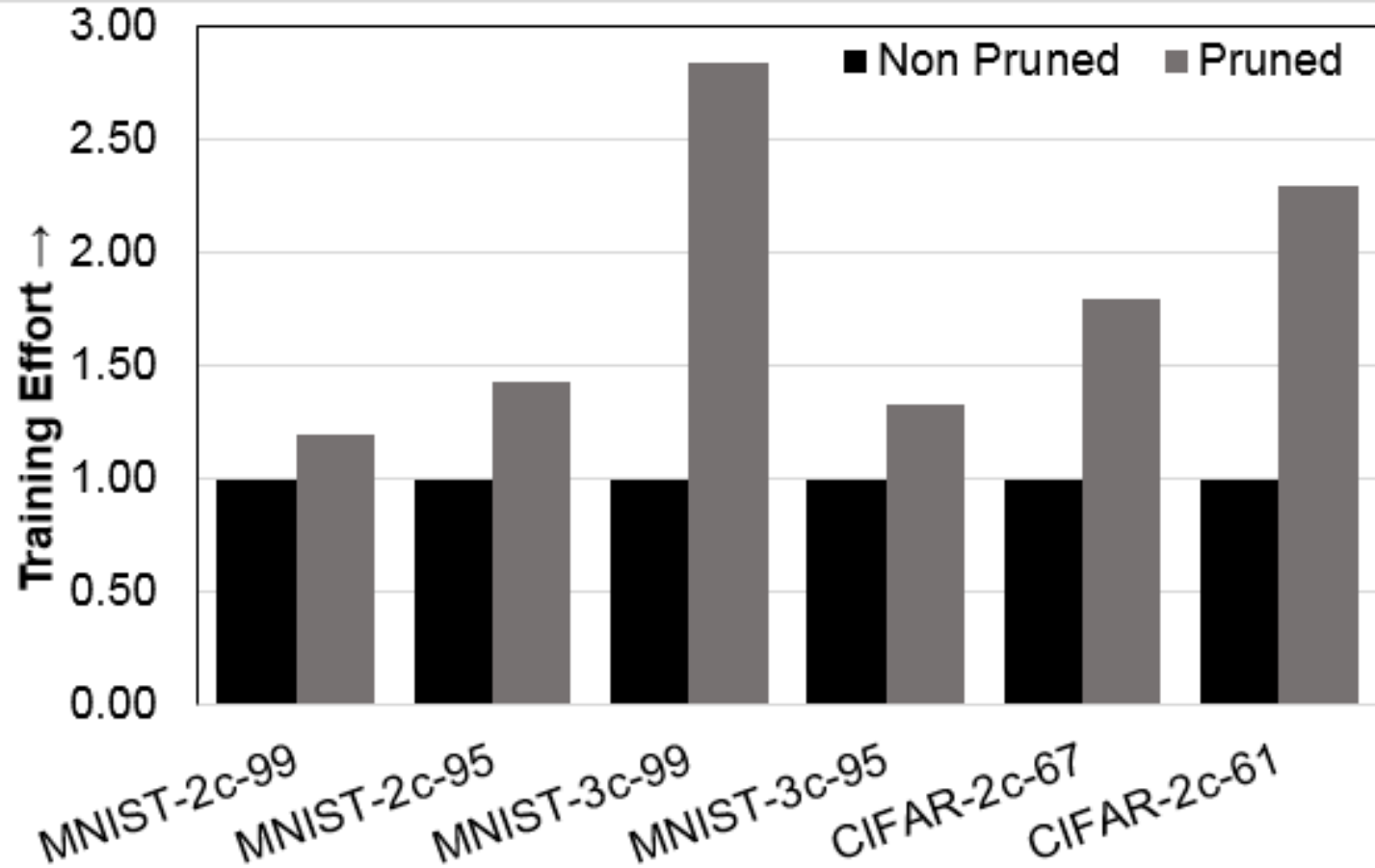
Result #1. Benefits from Pruning



Inferences:-

- Pruning can potentially lead to memory savings from 38% to 79% (61% on average) across all the benchmarks.
- Lesser the accuracy requirement, more is the scope of pruning.
- More number of layers (redundancy), more is the scope of pruning.

Result #2. Training Effort



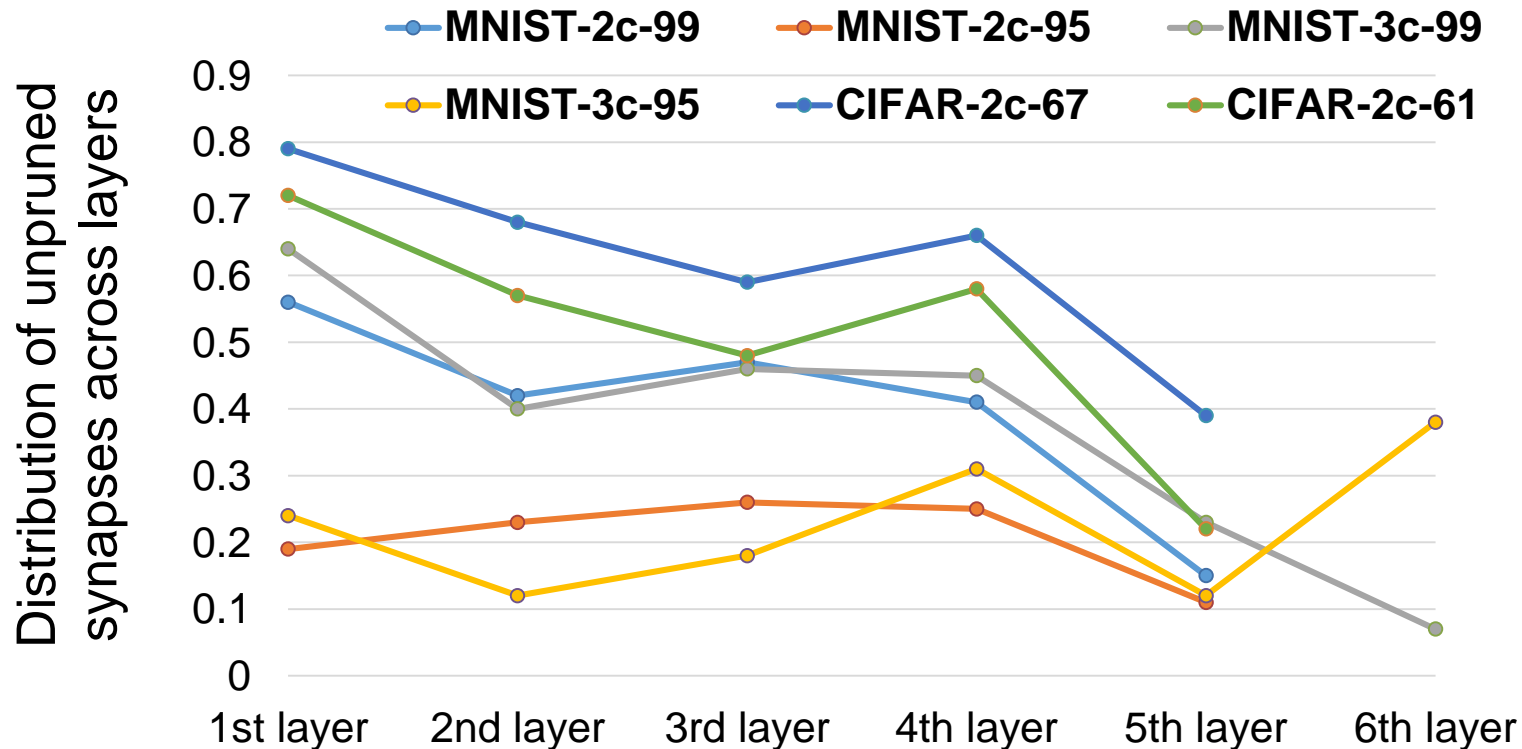
Pruning the neural network demands more training effort(time) to reach iso-accuracy.

➤ Training – Self-healing process

➤ Pruning – Approximation process



Result #3. Synapse distribution across layers after training



Key points:-

- This is one of the possible distributions which leads to our result and may not be the golden result.
- We do not infer the significance of one layer w.r.t. other layers in the NN from this graph as that requires further rigorous simulations.

Conclusions

- Pruning is a simple yet powerful technique to realize application aware flexible NN architectures (connectivity).
- Pruning and Training complement each other to obtain an optimally trained neural network.
- Amount of parameter reduction is a strong function of network structure, accuracy and training effort.

Future work

- Develop efficient pruning control techniques based on the training dynamics.
- Analyze pruning on bigger networks.

Questions ?