

### Visualization of Aligning Masks using Weight Symmetry



In Figure 1.a), a dense random init.,  $\mathbf{w}_{A}^{t=0}$ , converges to a dense solution,  $\mathbf{w}_{A}^{t=T}$ , which is then pruned with IMP resulting in the mask,  $\mathbf{m}_{A}$ . In Figure 1.b), permuting the mask,  $\pi(\mathbf{m}_A)$ , to match the (symmetric) basin in which the new initialization,  $\mathbf{w}_B^{t=0}$ , enables sparse training.

### Background

- . Lottery Ticket Hypothesis (LTH): identifies sparse sub-networks (i.e. binary masks) that, when trained independently, *can* match dense model performance. [1]
- 2. NNs are permutation invariant: swapping (i.e. permuting) neurons in a layer does not change the underlying function they compute.



In Figure 2., we show permutation symmetry in a single hidden layer NN, where the outputs, y and y' remain equivalent for the same input.

. Git Re-Basin claimed that NN loss landscapes nearly contain a single solution basin modulo permutations. [2]

### **Motivation**

Motivated by the goal of training a sparse model from a new random initialization.

The key limitations of LTH: a dense model must be first trained to get a mask, which can only be used with its original random initialization.

- LTH consistently converges to very similar solutions to the original pruned model, effectively relearning the same solution [4].

We seek to answer:

- generalization?
- different solutions?

he Lottery Ticket Hypothesis: Finding Sparse, Trainable Neural Networks, 2019 Ainsworth et al. Git re-basin: Merging models modulo permutation symmetries, 2023.

Paul et al. Unmasking the lottery ticket hypothesis: What's encoded in a winning ticket's mask?, 2023.

[4] Evci et al. Gradient flow in sparse neural networks and how lottery tickets win, 2022.



## Sparse Training from Random Initialization: Aligning Lottery Ticket Masks using Weight Symmetry $\textbf{Mohammed Adnan}^{*,1,3}, \textbf{Rohan Jain}^{*,1}, \textbf{Ekansh Sharma}^{2,3}, \textbf{Rahul G. Krishnan}^{2,3}, \textbf{Yani Ioannou}^1$

University of  $Calgary^1$  University of Toronto<sup>2</sup> Vector Institute<sup>3</sup>

### **Our Primary Findings**

"We found that LTH masks fail to generalize to new random initializations due to loss basin misalignment. To reuse an LTH mask with a different random initialization, we leverage permutation symmetries, to permute the mask to align with the new random initialization optimization basin."

- We find that a sparse model (with the permuted mask) with new random initialization can nearly match generalization performance of the LTH solution.
- We empirically demonstrate this on CIFAR-10/100 and ImageNet with VGG11 and ResNet models of varying widths.

### Methodology



### **Results**: ResNet50/ImageNet + VGG11/CIFAR-10

- **ResNet50**: permuted solution beats the naive solution across all sparsity levels, validating our hypothesis on large datasets.
- **VGG11**: increasing the rewind point, the permuted solution closely matches the accuracy of LTH, while naive solution significantly plateaus.



### **Results**: ResNet20/CIFAR-10 Observations

- Permuted solution outperforms the naive solution. Sparsity increasing: training becomes harder, widening the gap between permuted and naive solutions.
- Both the LTH & permuted solution do not perform well at a truly random init. (k = 0) but improves on increasing the rewind point until plateauing.
- Width increasing: the gap between training from random init. with the permuted mask & the LTH/ dense baseline decreases, unlike training with the naive mask.





In Figure 7.a): Various measures of function space similarity between the models. In Figure 7. b): 0-1 loss landscape of ResNet20x $\{4\}$ /CIFAR-100. In Figure 7.c): Error barrier between dense and LTH solutions after accounting for variance collapse.



# 

### **Results**: ResNet20/CIFAR-10 Plots



Larger width exhibits better linear mode connectivity (LMC). As the width of the model increases, the permutation matching algorithm gets more accurate, thereby reducing the loss barrier.



Although the mean test acc. of LTH is higher, ensemble of permuted models achieves better test acc. due to better functional diversity of permuted models.

We also show, modulo permutations, reusing the permuted mask leads to convergence in the same

• We show for a fixed init., the dense solution and corresponding LTH solution reside within the same loss basin when variance collapse is considered. This conclusion presents a new perspective compared to observations made by Paul et al. [3].

