



# Going Beyond Linear Mode Connectivity: The Layerwise Linear Feature Connectivity

**Zhanpeng Zhou**<sup>1</sup>, Yongyi Yang<sup>2</sup>, Xiaojiang Yang<sup>1</sup>, Junchi Yan<sup>1\*</sup>, Wei Hu<sup>2</sup>

<sup>1</sup>Shanghai Jiao Tong University, <sup>2</sup>University of Michigan



**SHANGHAI JIAO TONG  
UNIVERSITY**

**M UNIVERSITY OF MICHIGAN**

\*Corresponding author

# Background: Linear Mode Connectivity

## Linear Mode Connectivity (LMC)

Given dataset  $D$  and two modes  $\theta_A, \theta_B$  that  $\text{Err}_D(\theta_A) = \text{Err}_D(\theta_B)^*$ , two mode  $\theta_A$  and  $\theta_B$  satisfy the *linear mode connectivity* if

$$\forall \alpha \in [0, 1], \text{Err}_D(\alpha\theta_A + (1 - \alpha)\theta_B) \approx \text{Err}_D(\theta_A)$$

\* $\text{Err}_D(\theta)$  denotes the classification error of the network  $f(\theta; \cdot)$  on the dataset  $D$ .

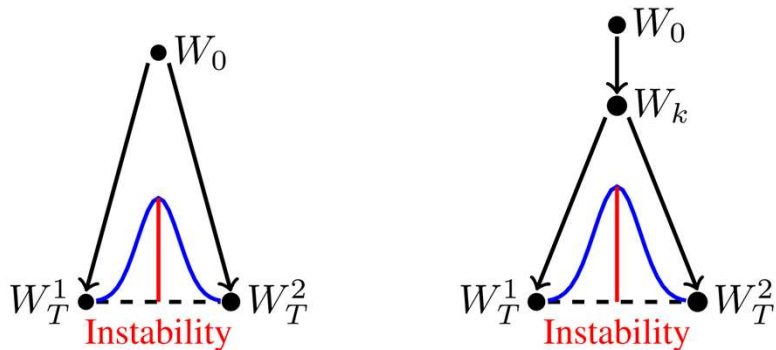


Fig. 1: Illustration of spawning method and LMC [1].

Frankle et al. [1] observed LMC for networks that are jointly trained for a short time before independent training (**spawning method**).

[1] Jonathan Frankle, Gintare Karolina Dziugaite, Daniel Roy, and Michael Carbin. Linear mode connectivity and the lottery ticket hypothesis.

# Background: Permutation Method

## Permutation Invariance.

Given an  $L$ -layer MLP  $f$ , we can permute the neurons of the MLP in each layer  $\ell \in [L]$  without changing its functionality ( $\pi = \{\mathbf{P}^{(\ell)}\}_{\ell \in [L]}$  are permutation matrices\*):

$$f(\boldsymbol{\theta}; \cdot) = f(\boldsymbol{\theta}'; \cdot), \text{ where } \boldsymbol{\theta} = \{\mathbf{W}^{(\ell)}\}_{\ell \in [L]}, \boldsymbol{\theta}' = \{\mathbf{W}'^{(\ell)}\}_{\ell \in [L]}$$

$$\forall \ell \in [L], \mathbf{W}'^{(\ell)} = \mathbf{P}^{(\ell)} \mathbf{W}^{(\ell)}, \mathbf{b}^{(\ell)} = \mathbf{P}^{(\ell)} \mathbf{b}^{(\ell)}, \mathbf{W}'^{(\ell+1)} = \mathbf{W}^{(\ell+1)} \mathbf{P}^{(\ell)}$$

\*Note that  $\mathbf{P}^{(0)}$  and  $\mathbf{P}^{(L)}$  are all fixed to be identity matrix.

Independently trained networks can be *linearly connected* when considering *permutation invariance* (**permutation methods**)[2, 3].

[2] Rahim Entezari, Hanie Sedghi, Olga Saukh, and Behnam Neyshabur. The role of permutation invariance in linear mode connectivity of neural networks.

[3] Samuel Ainsworth, Jonathan Hayase, and Siddhartha Srinivasa. Git re-basin: Merging models modulo permutation symmetries.

# Background: Permutation Method

Ainsworth et al. [3] proposed *weight matching* and *activation matching* to achieve LMC:

$$\text{weight matching}^*: \min_{\pi} \sum_{\ell=1}^L \left\| \mathbf{W}_A^{(\ell)} - \mathbf{P}^{(\ell)} \mathbf{W}_B^{(\ell)} \mathbf{P}^{(\ell-1)\top} \right\|_F^2$$

$$\text{Activation matching}^*: \min_{\pi} \sum_{\ell=1}^L \left\| \mathbf{H}_A^{(\ell)} - \mathbf{P}^{(\ell)} \mathbf{H}_B^{(\ell)} \right\|_F^2$$

\*We denote  $\ell$ -th layer feature as  $\mathbf{H}^{(\ell)}$  over the dataset  $D$ . Subscript  $\{A, B\}$  corresponds to modes  $\theta_A, \theta_B$ .

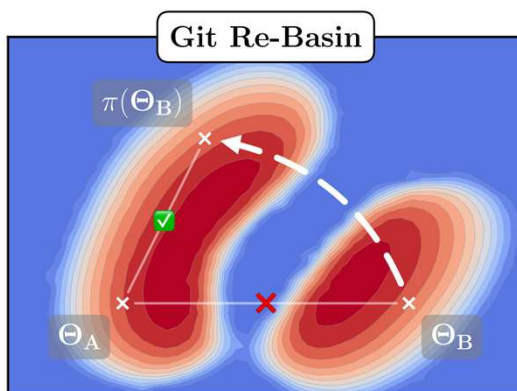
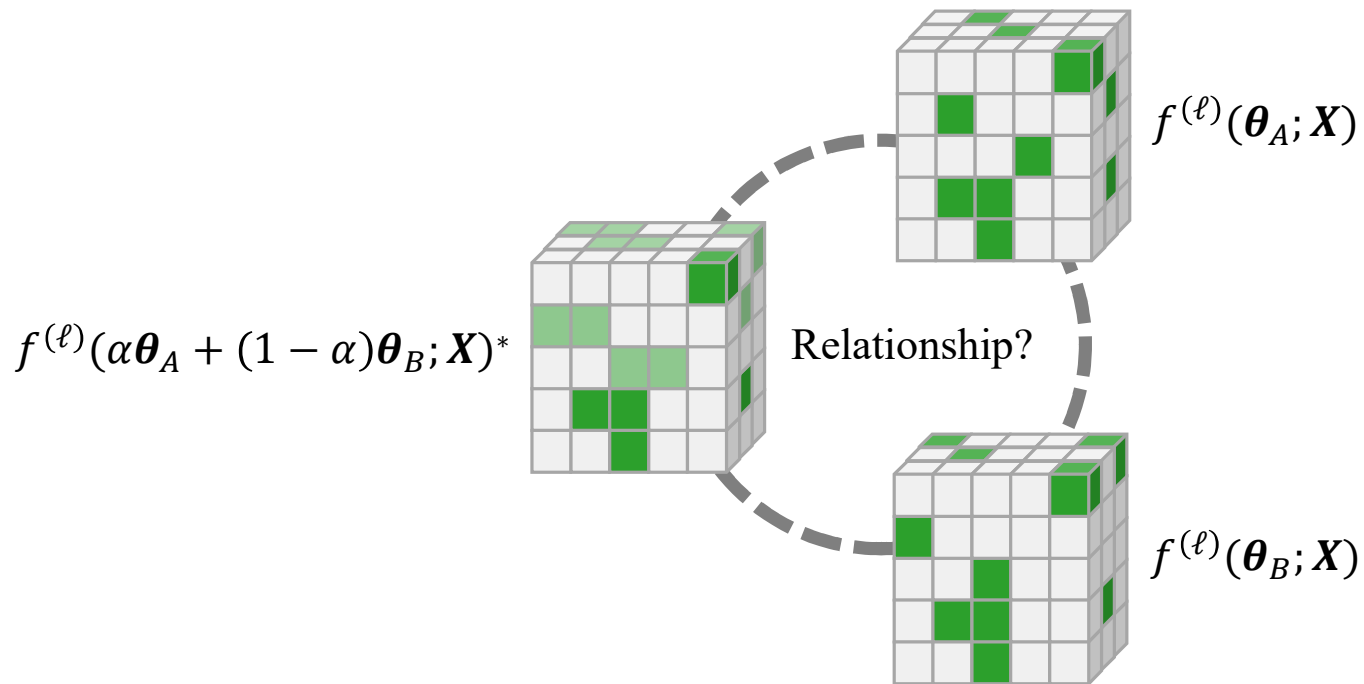


Fig. 2: Illustration of permutation [2].

# Motivation



*what happens to the internal features when we linearly interpolate the weights of two trained networks?*

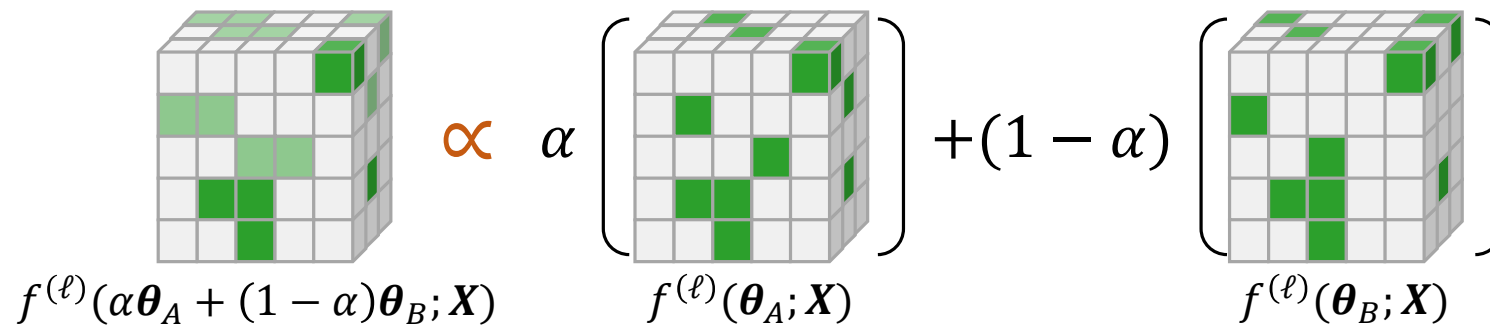
\* $f^{(\ell)}(\theta)$  denotes  $\ell$ -th layer feature of the network  $f(\theta; \cdot)$  over the dataset  $D$ .

# Layerwise Linear Feature Connectivity

## Layerwise Linear Feature Connectivity (LLFC)

Given dataset  $D$  and two modes  $\theta_A, \theta_B$  of an  $L$ -layer neural network  $f$ , the modes  $\theta_A$  and  $\theta_B$  are *layerwise linearly feature connected* if:

$$\forall \ell \in [L], \forall \alpha \in [0, 1], \exists c > 0, s. t., c f^{(\ell)}(\alpha \theta_A + (1 - \alpha) \theta_B) = \alpha f^{(\ell)}(\theta_A) + (1 - \alpha) f^{(\ell)}(\theta_B).$$



# Layerwise Linear Feature Connectivity

## LLFC always co-occurs with LMC in practice

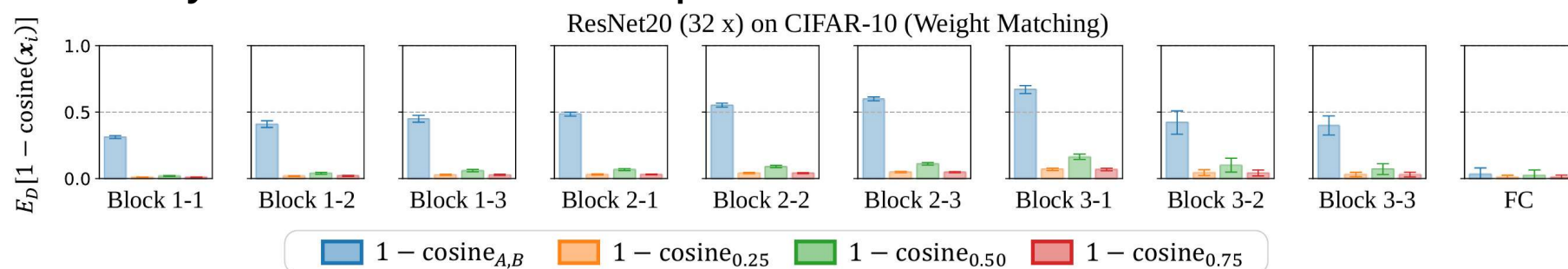


Fig. 3: Comparison of  $E_D[1 - \text{cosine}_\alpha(\mathbf{x}_i)]^*$  and  $E_D[1 - \text{cosine}_{A,B}(\mathbf{x}_i)]^*$ ,  $\alpha \in \{.25, .5, .75\}$ .

## Lemma (LLFC implies LMC)

Two modes  $\theta_A, \theta_B$  satisfy LLFC over dataset  $D$  and  $\max\{\text{Err}_D(\theta_A), \text{Err}_D(\theta_B)\} \leq \epsilon$

$$\forall \alpha \in [0, 1], \text{Err}_D(\alpha\theta_A + (1 - \alpha)\theta_B) \leq 2\epsilon.$$

\* $\text{cosine}_\alpha(\mathbf{x}_i) = \cos\langle f^{(\ell)}(\alpha\theta_A + (1 - \alpha)\theta_B; \mathbf{x}_i), \alpha f^{(\ell)}(\theta_A; \mathbf{x}_i) + (1 - \alpha)f^{(\ell)}(\theta_B; \mathbf{x}_i) \rangle$  and  $\text{cosine}_{A,B}(\mathbf{x}_i) = \cos\langle f^{(\ell)}(\theta_A; \mathbf{x}_i), f^{(\ell)}(\theta_B; \mathbf{x}_i) \rangle$

# Why LLFC Emerges?

Two simple conditions that leads to LLFC.

## Condition I: Weak Additivity for ReLU Activations

Given dataset  $D$ , the modes  $\theta_A$  and  $\theta_B$  satisfy *weak additivity for ReLU activations* if

$$\forall \ell \in [L], \forall \alpha \in [0,1], \sigma \left( \alpha \tilde{\mathbf{H}}_A^{(\ell)} + (1 - \alpha) \tilde{\mathbf{H}}_B^{(\ell)} \right) = \alpha \sigma \left( \tilde{\mathbf{H}}_A^{(\ell)} \right) + (1 - \alpha) \sigma \left( \tilde{\mathbf{H}}_B^{(\ell)} \right).^*$$

\*We denote  $\ell$ -th layer pre-activations as  $\tilde{\mathbf{H}}^{(\ell)}$  over the dataset  $D$  and ReLU activation as  $\sigma(\cdot)$ .

## Condition II: Commutativity

Given dataset  $D$ , the modes  $\theta_A$  and  $\theta_B$  satisfy *commutativity* if

$$\forall \ell \in [L], \mathbf{W}_A^{(\ell)} \mathbf{H}_A^{(\ell-1)} + \mathbf{W}_B^{(\ell)} \mathbf{H}_B^{(\ell-1)} = \mathbf{W}_B^{(\ell)} \mathbf{H}_A^{(\ell-1)} + \mathbf{W}_A^{(\ell)} \mathbf{H}_B^{(\ell-1)}.$$



# Why LLFC Emerges?

## Theorem (Condition I and II imply LLFC)

Given dataset  $D$ , if two modes  $\boldsymbol{\theta}_A$  and  $\boldsymbol{\theta}_B$  satisfy *weak additivity for ReLU activations* and *commutativity*, then

$$\forall \ell \in [L], \forall \alpha \in [0, 1], f^{(\ell)}(\alpha \boldsymbol{\theta}_A + (1 - \alpha) \boldsymbol{\theta}_B) = \alpha f^{(\ell)}(\boldsymbol{\theta}_A) + (1 - \alpha) f^{(\ell)}(\boldsymbol{\theta}_B).^*$$

*Weak additivity for ReLU activations* and *commutativity* are verified empirical for modes that satisfy LMC/LLFC.

# Justification of Permutation Method

Given a mode  $\boldsymbol{\theta}_A$  and a permuted mode  $\boldsymbol{\theta}'_B = \pi(\boldsymbol{\theta}_B)$  that satisfy LLFC, the *commutativity* is satisfied:

$$\forall \ell \in [L], \mathbf{W}_A^{(\ell)} \mathbf{H}_A^{(\ell-1)} + \mathbf{W}'_B^{(\ell)} \mathbf{H}'_B^{(\ell-1)} = \mathbf{W}'_B^{(\ell)} \mathbf{H}_A^{(\ell-1)} + \mathbf{W}_A^{(\ell)} \mathbf{H}'_B^{(\ell-1)} \quad (1)$$

Rewritten as:

$$\forall \ell \in [L], \left( \mathbf{W}_A^{(\ell)} - \mathbf{P}^{(\ell)} \mathbf{W}_B^{(\ell)} \mathbf{P}^{(\ell-1)\top} \right) \left( \mathbf{H}_A^{(\ell-1)} - \mathbf{P}^{(\ell)} \mathbf{H}_B^{(\ell-1)} \right) = 0 \quad (2)$$

## Connection to permutation methods

$$\text{weight matching: } \min_{\pi} \sum_{\ell=1}^L \left\| \mathbf{W}_A^{(\ell)} - \mathbf{P}^{(\ell)} \mathbf{W}_B^{(\ell)} \mathbf{P}^{(\ell-1)\top} \right\|_F^2$$

$$\text{Activation matching: } \min_{\pi} \sum_{\ell=1}^L \left\| \mathbf{H}_A^{(\ell)} - \mathbf{P}^{(\ell)} \mathbf{H}_B^{(\ell)} \right\|_F^2$$

The two objectives correspond to the two factors of above equation.

# Conclusion

## Conclusion

- Identify Layerwise Linear Feature Connectivity (LLFC)
- Investigate the underlying contributing factors to LLFC
- Obtain novel insights into permutation methods

## Future Directions

- Feature averaging methods
- Find a permutation directly enforcing the commutativity property
- [Going Beyond Neural Network Feature Similarity: The Network Feature Complexity and Its Interpretation Using Category Theory](#)