

Appendix of “On Linear Mode Connectivity of Mixture-of-Experts Architectures”

Table of Contents

A	On the Weight Spaces of Mixture-of-Experts Architecture	17
A.1	Weight Space of Mixture-of-Experts	17
A.2	Group Action on Weight Spaces of Mixture-of-Experts	18
B	Results on Mixture-of-Experts with Dense Gating	20
B.1	A Structural Property of Neural Networks with ReLU Activation	20
B.2	A Technical Lemma on Holomorphic Functions in \mathbb{C}^n	20
B.3	Functional Equivalence in Mixture-of-Experts with Dense Gating	21
C	Results on Mixture-of-Experts with Dense Gating	27
C.1	Strongly Distinctness Property	27
C.2	Functional Equivalence in Mixture-of-Experts with Sparse Gating	28
D	Technical Details for Sections 5 and 6	32
D.1	Proof for the Sufficiency of Permutation Invariance in LMC of MoE	32
D.2	Proof of the Permutation-Invariant Property for Equation (11)	33
D.3	Weight Matching Algorithm for Mixture-of-Experts	34
D.4	Formal formulation of DeepSeekMoE	34
E	Impact of Feedforward Reinitialization on Pretrained Transformer Performance	35
F	Experimental Details and Hyperparameters	37
G	Experimental Results	38
G.1	Verification of Linear Mode Connectivity across diverse configurations	38
G.1.1	Dense Mixture-of-Experts	40
G.1.2	Sparse Mixture-of-Experts	50
G.1.3	DeepSeek Mixture-of-Experts	57
G.2	Linear Mode Connectivity Analysis: Last Layer	64
G.3	Expert Matching Method	71
G.4	Ablation Study on Number of Layers	73
H	Broader Impact	76

622 A On the Weight Spaces of Mixture-of-Experts Architecture

623 Denote the input token dimension as d .

624 A.1 Weight Space of Mixture-of-Experts

625 **Expert functions.** We study expert functions $\mathcal{E} : \mathbb{R}^d \rightarrow \mathbb{R}^d$, which are realized by standard feedfor-
 626 ward neural networks employing the ReLU activation function. Specifically, each expert map $\mathcal{E}(\cdot; \theta)$
 627 is represented as a composition of affine transformations and ReLU nonlinearities:

$$\mathcal{E}(x; \theta) = f_L \circ \sigma \circ f_{L-1} \circ \cdots \circ \sigma \circ f_1(x), \quad (12)$$

628 where each affine layer is defined by

$$f_i(x) = xW_i + b_i, \quad \text{for } i = 1, \dots, L. \quad (13)$$

629 Here, σ denotes the ReLU activation function, applied component-wise, and $\theta = \{W_i, b_i\}_{i=1}^L$
 630 represents the full set of learnable parameters. The weight matrix $W_i \in \mathbb{R}^{n_{i-1} \times n_i}$ and bias vector
 631 $b_i \in \mathbb{R}^{n_i}$ parameterize the i -th layer, where $n_0 = n_L = d$ is the input and output dimensions,
 632 respectively. We have

$$\theta \in \prod_{i=1}^L (\mathbb{R}^{n_{i-1} \times n_i} \times \mathbb{R}^{n_i}) = \mathbb{R}^e, \quad \text{where } e = \sum_{i=1}^L (n_{i-1}n_i + n_i). \quad (14)$$

633 Define the *weight space of expert functions* as:

$$\Theta = \prod_{i=1}^L (\mathbb{R}^{n_{i-1} \times n_i} \times \mathbb{R}^{n_i}) = \mathbb{R}^e. \quad (15)$$

634 **Dense Mixture-of-Experts.** Given a positive integer n representing the number of experts, a
 635 *Mixture-of-Experts with Dense gating* (MoE) is defined as the function: $\mathcal{D} : \mathbb{R}^d \rightarrow \mathbb{R}^d$ defined by

$$\mathcal{D}(x; \{W_i, b_i, \theta_i\}_{i=1}^n) = \sum_{i=1}^n \text{softmax}_i(\{W_i x + b_i\}_{i=1}^n) \cdot \mathcal{E}(x; \theta_i). \quad (16)$$

636 The function \mathcal{D} is parameterized as $\mathcal{D}(x; \phi)$ where

$$\phi = \left((W_i, b_i), \theta_i \right)_{i=1, \dots, n} \in \left((\mathbb{R}^D \times \mathbb{R}) \times \Theta \right)^n. \quad (17)$$

637 Denote the *weight space of a Mixture of n Experts* as

$$\Phi(n) = \left((\mathbb{R}^D \times \mathbb{R}) \times \Theta \right)^n. \quad (18)$$

638 Varying the number of experts leads to a Mixture-of-Experts weight space that spans across expert
 639 sets of different sizes, denoted by

$$\Phi = \bigsqcup_{n=1}^{\infty} \Phi(n) = \bigsqcup_{n=1}^{\infty} \left((\mathbb{R}^D \times \mathbb{R}) \times \Theta \right)^n. \quad (19)$$

640 **Sparse Mixture-of-Experts.** Given a positive integer $k \leq n$, the Top- k map is defined by: for any
 641 vector $z = (z_1, \dots, z_n) \in \mathbb{R}^n$,

$$\text{Top-}k(z) = \{i_1, \dots, i_k\}, \quad (20)$$

642 where i_1, \dots, i_k are the indices corresponding to the k largest components of x . In the event of ties,
 643 we select smaller indices first. Using this, a *Mixture-of-Experts with Sparse gating* (SMoE) is the
 644 function $\mathcal{S} : \mathbb{R}^d \rightarrow \mathbb{R}^d$ defined by

$$\mathcal{S}(x; \{W_i, b_i, \theta_i\}_{i=1}^n) = \sum_{i \in T(x)} \text{softmax}_i(\{W_i x + b_i\}_{i \in T(x)}) \cdot \mathcal{E}(x; \theta_i). \quad (21)$$

645 where

$$T(x) = T(x; \{W_i, b_i\}_{i=1}^n) = \text{Top-}k\left((W_i x + b_i)_{i=1}^n\right). \quad (22)$$

646 The weight space for the Sparse Mixture-of-Experts is identical to that of the Dense Mixture-of-
 647 Experts, as the inclusion of the Top- k operator does not introduce any additional trainable parameters.
 648 However, the SMoE function \mathcal{S} is generally discontinuous, primarily due to the non-smooth nature
 649 of the Top- k selection—particularly in scenarios where multiple experts share identical gating scores
 650 (i.e., ties occur). To circumvent this issue and ensure well-defined behavior, we restrict our analysis
 651 to a subset of the input space \mathbb{R}^d where the top K gating scores are uniquely determined without
 652 ambiguity. In particular, for $\{W_i, b_i\}_{i=1}^n \in (\mathbb{R}^d \times \mathbb{R})^n$, we define the subset of \mathbb{R}^d such that the
 653 gating scores of all experts are pairwise distinct, i.e.,

$$\Omega(\{W_i, b_i\}_{i=1}^n) = \{x \in \mathbb{R}^d : (W_i x + b_i)_{i=1}^n \text{ are pairwise distinct}\}. \quad (23)$$

654 We provide one result characterizing the structure of this domain, and another result describing the
 655 behavior of the SMoE map when it is restricted to inputs within this domain.

656 **Proposition A.1.** *If $\{W_i, b_i\}$ are pairwise distinct for $i = 1, \dots, n$, then $\Omega(\{W_i, b_i\}_{i=1}^n)$ is an open
 657 and dense subset of \mathbb{R}^d .*

658 *Proof.* We have

$$\begin{aligned} \Omega(\{W_i, b_i\}_{i=1}^n) &= \{x \in \mathbb{R}^d : W_i x + b_i \text{ is pairwise distinct for all } i = 1, \dots, n\} \\ &= \bigcap_{1 \leq i < j \leq n} \{x \in \mathbb{R}^d : W_i x + b_i \neq W_j x + b_j\} \\ &= \bigcap_{1 \leq i < j \leq n} (\mathbb{R}^d \setminus \{x \in \mathbb{R}^d : W_i x + b_i = W_j x + b_j\}). \end{aligned} \quad (24)$$

659 Note that, the set

$$\{x \in \mathbb{R}^d : W_i x + b_i = W_j x + b_j\} = \{x \in \mathbb{R}^d : (W_i - W_j)x = b_j - b_i\}, \quad (25)$$

660 is either a hyperplane (when $W_i \neq W_j$) or empty (when $W_i = W_j$ but $b_i \neq b_j$). In both scenarios,
 661 its complement in \mathbb{R}^d forms an open and dense subset. By Equation (24), and using the fact that a
 662 finite intersection of open and dense subsets in \mathbb{R}^d remains open and dense, it follows that the set
 663 $\Omega(\{W_i, b_i\}_{i=1}^n)$ is itself open and dense in \mathbb{R}^d . \square

664 **Proposition A.2.** *If $\{W_i, b_i\}$ are pairwise distinct for $i = 1, \dots, n$, then the SMoE map \mathcal{S} , as given
 665 in Equation (21), is continuous over the domain $\Omega(\{W_i, b_i\}_{i=1}^n)$.*

666 *Proof.* Let $x \in \Omega(\{W_i, b_i\}_{i=1}^n)$. By the definition of this set, there exists an open neighborhood
 667 $U \subset \mathbb{R}^d$ containing x such that $U \subset \Omega(\{W_i, b_i\}_{i=1}^n)$ and

$$\text{Top-}k((W_i x + b_i)_{i=1}^n) = \text{Top-}k((W_i y + b_i)_{i=1}^n) \quad (26)$$

668 for all $y \in U$. This condition ensures that the sparse gating mechanism defined in Equa-
 669 tion (21) remains fixed within U , and consequently, the SMoE map is continuous on the domain
 670 $\Omega(\{W_i, b_i\}_{i=1}^n)$. \square

671 Propositions A.1 and A.2 serve as essential building blocks in the proof of Theorem C.4.

672 A.2 Group Action on Weight Spaces of Mixture-of-Experts

673 We define the group $G = G(n)$ as

$$G(n) = (\mathbb{R}^d \times \mathbb{R}) \times S_n, \quad (27)$$

674 which is the direct product of the additive groups \mathbb{R}^d and \mathbb{R} , and the symmetric group S_n on n
 675 elements. Each element $g \in G(n)$ can be written in the form

$$g = (c_W, c_b, \tau), \quad \text{where } c_W \in \mathbb{R}^d, \ c_b \in \mathbb{R}, \text{ and } \tau \in S_n. \quad (28)$$

676 The group $G(n)$ acts on the weight space $\Phi(n)$ as follows: for $g = (c_W, c_b, \tau) \in G(n)$ and $\phi \in \Phi(n)$
 677 given as in Equation (17), the action is defined by

$$g\phi := (W_{\tau(i)} + c_W, b_{\tau(i)} + c_b, \theta_{\tau(i)})_{i=1,\dots,n} \in \Phi(n). \quad (29)$$

678 This group action preserves the functionality of both the MoE and SMoE maps. The invariance arises
 679 from two key properties: the permutation invariance of the summation operator, and the translation
 680 invariance of the softmax function. We begin by establishing a result that characterizes the invariance
 681 of MoE maps under this group action.

682 **Proposition A.3** (Weight space invariance of Mixture-of-Experts). *The MoE map \mathcal{D} is $G(n)$ -*
 683 *invariance under the action of $G(n)$ on its weight space, i.e.,*

$$\mathcal{D}(\cdot; \phi) = \mathcal{D}(\cdot; g\phi). \quad (30)$$

684 *Proof.* Given $g = (c_W, c_b, \tau) \in G(n)$. For all $x \in \mathbb{R}^d$, we have

$$\begin{aligned} \mathcal{D}(x; g\theta) &= \sum_{i=1}^n \text{softmax}_i \left(\left\{ (W_{\tau(i)} + c_W) x + (b_{\tau(i)} + c_b) \right\}_{i=1}^n \right) \cdot \mathcal{E}(x; \theta_{\tau(i)}) \\ &= \sum_{i=1}^n \text{softmax}_i \left(\left\{ W_{\tau(i)} x + b_{\tau(i)} \right\}_{i=1}^n \right) \cdot \mathcal{E}(x; \theta_{\tau(i)}) \\ &= \sum_{i=1}^n \text{softmax}_i \left(\left\{ W_i x + b_i \right\}_{i=1}^n \right) \cdot \mathcal{E}(x; \theta_i) \\ &= \mathcal{D}(x; \theta). \end{aligned} \quad (31)$$

685 Thus, the proposition is proven. \square

686 The analysis of the SMoE architecture requires additional assumptions due to the inherent discontinu-
 687 ity of the Top- k selection operator. In what follows, we show that the SMoE map, when restricted to
 688 an appropriate subset of its domain, remains invariant under the group action of $G(n)$.

689 **Proposition A.4** (Weight space invariance of Sparse Mixture-of-Experts). *Given the SMoE map as*
 690 *defined in Equation (21), assume that the pairs $\{W_i, b_i\}$ are pairwise distinct for $i = 1, \dots, n$. Then,*
 691 *the set $\Omega(\{W_i, b_i\}_{i=1}^n)$ is invariant under the group action of $G(n)$, i.e.,*

$$\Omega(\{W_i, b_i\}_{i=1}^n) = \Omega(g\{W_i, b_i\}_{i=1}^n). \quad (32)$$

692 Moreover, the SMoE map, restricted to

$$\Omega(\{W_i, b_i\}_{i=1}^n), \quad (33)$$

693 is $G(n)$ -invariance under the action of $G(n)$ on its weight space, i.e.,

$$\mathcal{S}(\cdot; \theta) = \mathcal{S}(\cdot; g\theta) \quad \text{on} \quad \Omega(\{W_i, b_i\}_{i=1}^n). \quad (34)$$

694 *Proof.* Let $g = (c_W, c_b, \tau) \in G(n)$. We first verify that the group action preserves the set
 695 $\Omega(\{W_i, b_i\}_{i=1}^n)$. Indeed,

$$\begin{aligned} &\Omega(\{W_i, b_i\}_{i=1}^n) \\ &= \{x \in \mathbb{R}^d : W_i x + b_i \text{ is pairwise distinct for all } i = 1, \dots, n\} \\ &= \{x \in \mathbb{R}^d : (W_{\tau(i)} + c_W) x + (b_{\tau(i)} + c_b) \text{ is pairwise distinct for all } i = 1, \dots, n\} \\ &= \Omega(g\{W_i, b_i\}_{i=1}^n). \end{aligned} \quad (35)$$

696 Let

$$gT(x) = \text{Top-}k \left(\left((W_{\tau(i)} + c_W) x + (b_{\tau(i)} + c_b) \right)_{i=1}^n \right). \quad (36)$$

697 For all $x \in \Omega(\{W_i, b_i\}_{i=1}^n)$, we have $gT(x) = \tau(T(x))$. The desired invariance result for the SMoE
 698 map then follows by applying the same reasoning as in Proposition A.3. \square

Remark A.5. While the group action on Mixture-of-Experts (MoE) models is formally introduced in Equation (22), it does not fully capture all symmetries inherent to these architectures. In particular, each expert network possesses internal neuron permutations that preserve its functional behavior—a well-documented property in the literature on neural networks [7, 23, 29, 9, 60, 11, 62, 13, 77, 82, 37, 15, 27, 50, 5, 6]. Nonetheless, since the distinguishing feature of MoE models is their input-dependent gating function, our focus centers on the symmetries associated with the gate itself, treating expert-level invariances as part of a well-established theoretical foundation.

B Results on Mixture-of-Experts with Dense Gating

B.1 A Structural Property of Neural Networks with ReLU Activation

A *polytope* is a geometric entity bounded by flat surfaces, which can be either finite (bounded) or infinite (unbounded) in extent. We define the notion of *local affineness* as follows: a function $f : \mathbb{R}^d \rightarrow \mathbb{R}^D$ is said to be locally affine if there exists a partition of \mathbb{R}^d into a finite set of polytopes such that on each polytope, the function f agrees with an affine transformation from \mathbb{R}^d to \mathbb{R}^D .

Remark B.1. While the term *local affineness* may take on different interpretations in other contexts, its usage here is unambiguous within the scope of this work.

We examine the local affineness property in ReLU neural networks. Consider a feedforward neural network $f : \mathbb{R}^{n_0} \rightarrow \mathbb{R}^{n_L}$ constructed from affine maps interleaved with ReLU activations, given by

$$f = f_L \circ \sigma \circ f_{L-1} \circ \cdots \circ \sigma \circ f_1,$$

where each map $f_i : \mathbb{R}^{d_{i-1}} \rightarrow \mathbb{R}^{d_i}$ is affine and has the form $f_i(x) = W_i x + b_i$, and σ denotes the ReLU function applied elementwise.

The combination of these affine transformations and ReLU activations partitions the domain \mathbb{R}^{d_0} into finitely many convex polytopes. On each such region, the activation status of the ReLU units—whether they transmit their input or suppress it to zero—remains fixed. Consequently, the network reduces to a simpler form in which each ReLU behaves either as the identity or the zero function. Given that both ReLU and affine maps are piecewise linear and closed under composition, the overall function f is affine within each activation region.

Hence, the network satisfies the following property:

$$f(x) = A_i x + b_i, \quad \text{for all } x \in P_i, \quad (37)$$

where the domain is partitioned into polytopes $\{P_i\}_{i=1}^m$, and A_i, b_i define the affine transformation active within region P_i .

Let ∂P_i denote the boundary of the polytope P_i . Then the set

$$\mathbb{R}^{d_0} \setminus \bigcup_{i=1}^m \partial P_i \quad (38)$$

is open and dense in \mathbb{R}^{d_0} ; in other words, the union of the interiors of the polytopes covers a set that is both open and dense.

Now consider a finite collection of ReLU networks $\{f^{(k)}\}_{k=1}^n$. Since the intersection of finitely many open dense subsets is itself open and dense, there exists a subset $\Omega \subset \mathbb{R}^{d_0}$ such that, for every point $x \in \Omega$, there is a neighborhood around x on which all functions $f^{(k)}$ act as affine maps.

B.2 A Technical Lemma on Holomorphic Functions in \mathbb{C}^n

A function $f : \mathbb{C}^n \rightarrow \mathbb{C}$ is said to be *holomorphic* on \mathbb{C}^n if it is complex differentiable at every point in \mathbb{C}^n . A function is called *meromorphic* on \mathbb{C}^n if it can be locally written as a quotient of two holomorphic functions, where the denominator is not identically zero. The collection of all holomorphic functions on \mathbb{C}^n forms an integral domain, denoted by \mathcal{H} , while the set of meromorphic functions on \mathbb{C}^n forms a field, denoted by \mathcal{F} . In particular, \mathcal{F} is the field of fractions of the integral domain \mathcal{H} . Let $\mathbb{C}[x] = \mathbb{C}[x_1, \dots, x_n]$ denote the polynomial ring in n complex variables, and let $\mathbb{C}(x) = \mathbb{C}(x_1, \dots, x_n)$ denote the corresponding field of rational functions. Then $\mathbb{C}[x]$ is a subring of \mathcal{H} and remains an integral domain, while $\mathbb{C}(x)$ is a subfield of \mathcal{F} and represents the field of fractions of $\mathbb{C}[x]$.

Remark B.2. For any polynomial $p \in \mathbb{C}[x]$, the exponential e^p defines a holomorphic function on \mathbb{C}^n , i.e. $e^p \in \mathcal{D}$.

Since $\mathbb{C}(x) \subset \mathcal{F}$, we can view \mathcal{F} as a vector space over $\mathbb{C}(x)$. The following lemma addresses the linear independence of exponential functions of polynomials in this vector space setting.

Lemma B.3. Let p_1, \dots, p_N be polynomials in $\mathbb{C}[x]$ such that $p_i - p_j$ is nonconstant whenever $i \neq j$. Then the functions e^{p_1}, \dots, e^{p_N} , viewed as elements of \mathcal{F} , are linearly independent over $\mathbb{C}(x)$.

Proof. We proceed by induction on N . The base case $N = 1$ is immediate, since $e^p \neq 0$ for any polynomial p . Assume the result holds for any smaller number of exponentials. Consider polynomials $r_1, \dots, r_N \in \mathbb{C}[x]$ such that

$$r_1 \cdot e^{p_1} + \dots + r_N \cdot e^{p_N} = 0. \quad (39)$$

We aim to prove that all $r_i = 0$. Suppose not; then at least one r_i is nonzero. Without loss of generality, assume $r_N \neq 0$. Dividing through by $r_N \cdot e^{p_N}$ gives

$$\frac{r_1}{r_N} \cdot e^{p_1 - p_N} + \dots + \frac{r_{N-1}}{r_N} \cdot e^{p_{N-1} - p_N} + 1 = 0. \quad (40)$$

Differentiating both sides with respect to x_i for each $i = 1, \dots, n$, we obtain:

$$\sum_{j=1}^{N-1} \left(\frac{\partial}{\partial x_i} \left(\frac{r_j}{r_N} \right) + \frac{r_j}{r_N} \cdot \frac{\partial}{\partial x_i} (p_j - p_N) \right) \cdot e^{p_j - p_N} = 0. \quad (41)$$

Note that each term

$$\frac{\partial}{\partial x_i} \left(\frac{r_j}{r_N} \right) + \frac{r_j}{r_N} \cdot \frac{\partial}{\partial x_i} (p_j - p_N) \quad (42)$$

belongs to $\mathbb{C}(x)$. Now, the polynomials $p_1 - p_N, \dots, p_{N-1} - p_N$ are pairwise distinct and nonconstant. Thus, by the induction hypothesis, the exponentials $e^{p_j - p_N}$ are linearly independent over $\mathbb{C}(x)$ for $j = 1, \dots, N-1$. Therefore, from equation (41), we must have

$$\frac{\partial}{\partial x_i} \left(\frac{r_j}{r_N} \right) + \frac{r_j}{r_N} \cdot \frac{\partial}{\partial x_i} (p_j - p_N) = 0, \quad (43)$$

for all $i = 1, \dots, n$ and $j = 1, \dots, N-1$, which implies

$$\frac{\partial}{\partial x_i} \left(\frac{r_j}{r_N} \cdot e^{p_j - p_N} \right) = 0. \quad (44)$$

Hence, for each $j = 1, \dots, N-1$, the function

$$\frac{r_j}{r_N} \cdot e^{p_j - p_N} = c_j \in \mathbb{C} \quad (45)$$

must be constant. If $c_j \neq 0$, then both $r_j \neq 0$ and $e^{p_j - p_N} = \frac{c_j r_N}{r_j}$ must hold. But this forces $e^{p_j - p_N}$ to be constant, which contradicts the assumption that $p_j - p_N$ is nonconstant. Thus, $c_j = 0$, which implies $r_j = 0$ for all $j = 1, \dots, N-1$. Plugging back into equation (40) gives a contradiction, as $1 \neq 0$. Therefore, all $r_i = 0$, completing the proof. \square

Remark B.4. This lemma plays a key role and will be applied repeatedly in the proofs of Theorem B.5 and Theorem C.4.

B.3 Functional Equivalence in Mixture-of-Experts with Dense Gating

The following result establishes the equivalence between two sets of weights that define the same MoE map. Certain assumptions are introduced for technical reasons, and their justification is provided in Remark B.6.

Theorem B.5 (Functional equivalence in Mixture-of-Experts with Dense Gating). Suppose ϕ, ϕ' define the same MoE function, i.e. $\mathcal{D}(\cdot; \phi) = \mathcal{D}(\cdot; \phi')$. Assume that the following conditions hold:

1. Both $\{\mathcal{E}(\cdot; \theta_i)\}_{i=1}^n$ and $\{\mathcal{E}(\cdot; \theta'_i)\}_{i=1}^{n'}$ consist of pairwise distinct functions;

774 2. Both $\{W_i - W_j\}_{1 \leq i < j \leq n}$ and $\{W'_i - W'_j\}_{1 \leq i < j \leq n'}$ consist of pairwise distinct vector of \mathbb{R}^d .

775 Then $n = n'$, and there exists $g = (c_W, c_b, \tau) \in G(n)$ such that for all $i = 1, \dots, n$, we have

776 $W'_i = W_{\tau(i)} + c_W$, $b'_i = b_{\tau(i)} + c_b$, and $\mathcal{E}(\cdot; \theta'_i) = \mathcal{E}(\cdot; \theta_{\tau(i)})$ on \mathbb{R}^d .

777 *Proof.* To enhance clarity, we begin by outlining the main steps of the proof at a high level:

- 778 1. We first expand the equality $\mathcal{D}(\cdot; \phi) = \mathcal{D}(\cdot; \phi')$ and introduce simplified notation to stream-
779 line the subsequent derivations.
- 780 2. We then observe that each expert function can be locally characterized as an affine map.
- 781 3. Next, we prove that $n = n'$ and establish the existence of a permutation τ and a transforma-
782 tion c_W with the required properties.
- 783 4. We proceed to verify the equivalence between the two sets of experts.
- 784 5. Finally, we show that a transformation c_b satisfying the necessary conditions can be con-
785 structed.

786 We now proceed with the detailed derivations and justifications corresponding to each of the five
787 steps above.

788 **Step 1.** Given that $\mathcal{D}(\cdot; \phi) = \mathcal{D}(\cdot; \phi')$, we have

$$\sum_{i=1}^n \text{softmax}_i(\{W_i x + b_i\}_{i=1}^n) \cdot \mathcal{E}(x; \theta_i) = \sum_{i=1}^{n'} \text{softmax}_i(\{W'_i x + b'_i\}_{i=1}^{n'}) \cdot \mathcal{E}(x; \theta'_i), \quad (46)$$

789 for all $x \in \mathbb{R}^d$. Define

$$\mathcal{E}_i(\cdot) = \mathcal{E}(\cdot; \theta_i), \quad \text{and} \quad \mathcal{E}'_i(\cdot) = \mathcal{E}(\cdot; \theta'_i). \quad (47)$$

790 By expanding the softmax terms in Equation (46), we obtain

$$\sum_{i=1}^n \frac{e^{W_i x + b_i}}{\sum_{j=1}^n e^{W_j x + b_j}} \cdot \mathcal{E}_i(x) = \sum_{i=1}^{n'} \frac{e^{W'_i x + b'_i}}{\sum_{j=1}^{n'} e^{W'_j x + b'_j}} \cdot \mathcal{E}'_i(x). \quad (48)$$

791 Multiplying both sides by the respective denominators yields

$$\left(\sum_{j=1}^{n'} e^{W'_j x + b'_j} \right) \cdot \left(\sum_{i=1}^n e^{W_i x + b_i} \cdot \mathcal{E}_i(x) \right) = \left(\sum_{j=1}^n e^{W_j x + b_j} \right) \cdot \left(\sum_{i=1}^{n'} e^{W'_i x + b'_i} \cdot \mathcal{E}'_i(x) \right), \quad (49)$$

792 which can be rewritten as

$$\sum_{i=1}^n \sum_{j=1}^{n'} e^{(W_i + W'_j)x + (b_i + b'_j)} \cdot (\mathcal{E}_i(x) - \mathcal{E}'_j(x)) = 0. \quad (50)$$

793 **Step 2.** Since each function \mathcal{E}_i and \mathcal{E}'_j is locally affine, it follows from the result in Appendix B.1
794 that there exists a dense open subset $\Omega \subset \mathbb{R}^d$ such that: for any point $a \in \Omega$, one can find an open
795 neighborhood $U \subset \Omega$ containing a on which all functions \mathcal{E}_i and \mathcal{E}'_j are affine. Consequently, each of
796 these functions agrees with a polynomial on U . That is, there exists a family of open sets $\{U_k\}_{k \in I}$
797 covering Ω , so that

$$\Omega = \bigcup_{k \in I} U_k, \quad (51)$$

798 and for each set $U = U_k$ in this cover, there exist polynomials $p_{U,i}, p'_{U,j} \in \mathbb{R}[x]$ such that

$$\mathcal{E}_i(x) = p_{U,i}(x), \quad \text{and} \quad \mathcal{E}'_j(x) = p'_{U,j}(x) \quad \text{for all } x \in U. \quad (52)$$

799 Substituting into Equation (50) yields:

$$\sum_{i=1}^n \sum_{j=1}^{n'} e^{(W_i + W'_j)x + (b_i + b'_j)} \cdot (p_{U,i}(x) - p'_{U,j}(x)) = 0 \quad \text{for all } x \in U. \quad (53)$$

800 Observe that the expression on the left-hand side above defines a holomorphic function. By the
801 Identity Theorem for Holomorphic Functions (see [1, 66, 16, 74]), we conclude that:

$$\sum_{i=1}^n \sum_{j=1}^{n'} e^{(W_i + W'_j)x + (b_i + b'_j)} \cdot (p_{U,i}(x) - p'_{U,j}(x)) = 0 \quad \text{for all } x \in \mathbb{C}^d. \quad (54)$$

802 **Step 3.** According to Assumption 2, the sets $\{W_i\}_{i=1}^n$ and $\{W'_j\}_{j=1}^{n'}$ are composed of mutually
803 distinct elements. As a result, there exists a direction

$$\alpha \in \mathbb{S}^{d-1} = \{x \in \mathbb{R}^d : \|x\|_2 = 1\}, \quad (55)$$

804 such that the projected values $\{W_i \alpha\}_{i=1}^n$ and $\{W'_j \alpha\}_{j=1}^{n'}$ are comprised of n and n' distinct real
805 numbers, respectively. Without loss of generality, we may relabel the indices so that:

$$W_1 \alpha < W_2 \alpha < \dots < W_n \alpha \quad \text{and} \quad W'_1 \alpha < W'_2 \alpha < \dots < W'_{n'} \alpha. \quad (56)$$

806 Furthermore, observe that the problem setting and the preceding equations are invariant under
807 translations of the form $W'_j \mapsto W'_j + c_W$ for a constant vector $c_W \in \mathbb{R}^d$. Hence, we can assume
808 without loss of generality that $W_1 = W'_1$. Under this assumption, we aim to demonstrate that $n = n'$
809 and $W_i = W'_i$ for all $i = 1, \dots, n$. Toward that goal, we begin by showing that $W_i = W'_i$ for each
810 $i = 1, \dots, \min\{n, n'\}$ using mathematical induction.

811 *Base case.* By assumption, we have $W_1 = W'_1$, so the base case holds trivially.

812 *Auxiliary step for induction.* For every index pair $(i, j) \neq (1, 1)$, we have the inequality:

$$W_1 \alpha + W'_1 \alpha < W_i \alpha + W'_j \alpha. \quad (57)$$

813 Thus, the sum $W_1 + W'_1$ differs from $W_i + W'_j$ for all $(i, j) \neq (1, 1)$. Applying Equation (54) in
814 conjunction with Lemma B.3, we conclude that

$$p_{U,1} = p'_{U,1}. \quad (58)$$

815 *Inductive step.* Assume that $W_i = W'_i$ holds for all $1 \leq i < k$, where k is an integer such that
816 $1 < k \leq \min\{n, n'\}$. Suppose, for contradiction, that $W_k \neq W'_k$. We analyze the terms $W_1 + W'_k$
817 and $W_k + W'_1$. Under our assumption, these two quantities must differ. Without loss of generality,
818 we may assume that

$$W_1 \alpha + W'_k \alpha \leq W_k \alpha + W'_1 \alpha. \quad (59)$$

819 • For all index pairs (i, j) with $i \geq k$, we have

$$W_1 \alpha + W'_k \alpha \leq W_k \alpha + W'_1 \alpha \leq W_i \alpha + W'_j \alpha. \quad (60)$$

820 Equality holds if and only if $(i, j) = (k, 1)$. Moreover, since $W_1 + W'_k$ and $W_k + W'_1$ are
821 distinct, it follows that $W_1 + W'_k$ differs from $W_i + W'_j$ for all (i, j) with $i \geq k$.

822 • For all (i, j) with $j \geq k$, we have

$$W_1 \alpha + W'_k \alpha \leq W_i \alpha + W'_j \alpha. \quad (61)$$

823 Equality occurs only when $(i, j) = (1, k)$. Therefore, $W_1 + W'_k$ is distinct from $W_i + W'_j$
824 for all $(i, j) \neq (1, k)$ with $j \geq k$.

825 • For all (i, j) such that $i, j < k$, we claim that $W_1 + W'_k$ does not equal $W_i + W'_j$. Suppose,
826 for contradiction, that

$$W_1 + W'_k = W_i + W'_j \quad (62)$$

827 for some pair (i, j) with $i, j < k$. Then by the induction hypothesis,

$$W'_1 + W'_k = W'_i + W'_j. \quad (63)$$

828 Rearranging terms, we obtain

$$W'_1 - W'_j = W'_i - W'_k, \quad (64)$$

829 which contradicts the assumption that all such differences are pairwise distinct, since
830 $(1, j) \neq (i, k)$.

831 These observations collectively show that $W_1 + W'_k$ differs from $W_i + W'_j$ for all $(i, j) \neq (1, k)$.
832 Applying Equation (54) together with Lemma B.3, we conclude that

$$p_{U,1} = p'_{U,k}. \quad (65)$$

833 Additionally, from Equation (58), we already have

$$p'_{U,1} = p'_{U,k}. \quad (66)$$

834 This implies that $\mathcal{E}'_1 = \mathcal{E}'_k$ on U . Since this holds for every open set U in the covering $\{U_k\}_{k \in I}$, it
835 follows that $\mathcal{E}'_1 = \mathcal{E}'_k$ on Ω . Because Ω is dense in \mathbb{R}^d and each \mathcal{E}'_j is continuous, we conclude that
836 $\mathcal{E}'_1 = \mathcal{E}'_k$ on all of \mathbb{R}^d . This contradicts the assumption that the \mathcal{E}'_j functions are pairwise distinct.
837 Therefore, our initial assumption must be false, and we conclude that $W_1 + W'_k = W_k + W'_1$, which
838 implies $W_k = W'_k$.

839 *Conclusion.* By induction, we have established that $W_i = W'_i$ for all $i = 1, \dots, \min\{n, n'\}$. It
840 remains to show that $n = n'$. Suppose, for contradiction, that $n < n'$. Consider the sum $W_1 + W'_{n'}$.
841 We claim that this quantity is distinct from every $W_i + W'_j$ with $(i, j) \neq (1, n')$. Assume otherwise,
842 that

$$W_1 + W'_{n'} = W_i + W'_j \quad (67)$$

843 for some pair $(i, j) \neq (1, n')$. Using the inductive assumption that $W_i = W'_i$ for all $i \leq n$, we
844 deduce

$$W'_1 + W'_{n'} = W'_i + W'_j, \quad (68)$$

845 which implies

$$W'_1 - W'_j = W'_i - W'_{n'}. \quad (69)$$

846 This leads to a contradiction, as it violates the assumption that the differences $W'_i - W'_j$ are pairwise
847 distinct. Hence, $W_1 + W'_{n'}$ must be distinct from all other sums $W_i + W'_j$ where $(i, j) \neq (1, n')$.
848 Then, by Equation (54) and Lemma B.3, it follows that

$$p_{U,1} = p'_{U,n'}. \quad (70)$$

849 In addition, from Equation (58), we already know that

$$p'_{U,1} = p'_{U,n'}. \quad (71)$$

850 Consequently, we obtain $\mathcal{E}'_1 = \mathcal{E}'_{n'}$ on U . Since this holds on every open set U in the covering
851 $\{U_k\}_{k \in I}$, it extends to Ω , and by continuity, to all of \mathbb{R}^d . This contradicts the assumption that the
852 expert functions \mathcal{E}'_j are pairwise distinct. Therefore, our assumption must be false, and we conclude
853 that $n = n'$. Finally, the reordering of indices and the translation applied to the set $\{W'_j\}_{j=1}^{n'}$
854 throughout the argument confirm the existence of a permutation $\tau \in S_n$ and a translation vector
855 $c_W \in \mathbb{R}^d$.

856 **Step 4.** We now establish that $\mathcal{E}_i = \mathcal{E}'_i$ on \mathbb{R}^d for each $i = 1, \dots, n$. From **Step 3**, we already
857 have that $n = n'$ and $W_i = W'_i$ for all indices in this range. Consider any pair (i, j) . If it holds
858 that $W_i + W'_j = W_{i'} + W'_{j'}$, then (i', j') must be either (i, j) or (j, i) . In particular, this implies
859 that $W_i + W'_i$ is distinct from $W_j + W'_k$ for all $(j, k) \neq (i, i)$. Invoking Equation (54) together with
860 Lemma B.3, we obtain

$$p_{U,i} = p'_{U,i}. \quad (72)$$

861 This argument parallels that used in **Step 3**, and applying the same reasoning, we conclude that
862 $\mathcal{E}_i = \mathcal{E}'_i$ on \mathbb{R}^d . Since this holds for all $i = 1, \dots, n$, the result follows.

863 **Step 5.** It remains to prove the existence of a constant $c_b \in \mathbb{R}$ such that

$$b'_i = b_i + c_b \quad \text{for all } i = 1, \dots, n. \quad (73)$$

864 Recall from **Step 4** that if $W_i + W'_j = W_{i'} + W'_{j'}$, then it must be that $(i', j') = (i, j)$ or (j, i) . Using
865 this structural property, along with Equation (50), Lemma B.3, and the identity $\mathcal{E}_i = \mathcal{E}'_i$ established
866 in **Step 4**, we derive the following equality:

$$e^{(W_i + W'_j)x + (b_i + b'_j)} \cdot (\mathcal{E}_i(x) - \mathcal{E}_j(x)) + e^{(W_j + W'_i)x + (b_j + b'_i)} \cdot (\mathcal{E}_j(x) - \mathcal{E}_i(x)) = 0, \quad (74)$$

867 valid for all index pairs (i, j) . Since $\mathcal{E}_i \neq \mathcal{E}_j$ whenever $i \neq j$, there exists a point $x_0 \in \mathbb{R}^d$ for which
868 $\mathcal{E}_i(x_0) \neq \mathcal{E}_j(x_0)$. Plugging $x = x_0$ into Equation (74) and simplifying by factoring out the common
869 nonzero difference, we find:

$$e^{b_i + b'_j} = e^{b_j + b'_i}, \quad (75)$$

870 which implies

$$b_i + b'_j = b_j + b'_i, \quad (76)$$

871 and hence

$$b_i - b'_i = b_j - b'_j. \quad (77)$$

872 This shows that the difference $b_i - b'_i$ remains constant for all i . Letting $c_b := b'_1 - b_1$, we obtain

$$b'_i = b_i + c_b \quad \text{for all } i = 1, \dots, n. \quad (78)$$

873 This completes the proof of Theorem B.5.

874 □

875 **Remark B.6** (Rationale behind the assumptions in Theorem B.5). In modeling architectures, it is
876 important that the symmetry group arises from the structure of the model itself rather than from
877 specific, possibly degenerate, parameter choices. That is, the symmetry group should act globally
878 and consistently across the entire weight space. This requirement motivates the following conditions
879 in Theorem B.5:

- 880 1. Both $\{\mathcal{E}(\cdot; \theta_i)\}_{i=1}^n$ and $\{\mathcal{E}(\cdot; \theta'_i)\}_{i=1}^{n'}$ consist of pairwise distinct functions;
- 881 2. Both $\{W_i - W_j\}_{1 \leq i < j \leq n}$ and $\{W'_i - W'_j\}_{1 \leq i < j \leq n'}$ consist of pairwise distinct vector of \mathbb{R}^d .

882 We now elaborate on the purpose of these assumptions.

883 *Assumption 1.* If this condition is violated—for instance, if two experts compute the same function and
884 are assigned identical gating values—then permuting these experts leaves the model output unchanged.
885 Such permutations introduce additional, non-essential elements into the symmetry group, which
886 we refer to as spurious symmetries. These do not correspond to genuine structural invariances but
887 rather arise from degenerate parameter settings that represent singular points in the model's parameter
888 space.

889 *Assumption 2.* Assumption 2 addresses a more nuanced issue: it rules out cases where linear
890 dependencies among gating weight vectors may cause multiple experts to exhibit indistinguishable
891 gating behavior. Although such situations may not be as immediately intuitive as those excluded by
892 Assumption 1, they too can artificially enlarge the symmetry group beyond its intended form. To
893 illustrate this, we present an explicit example. Let $d = 1$ and $n = n' = 3$, and consider parameter
894 configurations ϕ and ϕ' such that:

- 895 • $W_1 = W'_1 = -1$, $W_2 = W'_2 = 0$, and $W_3 = W'_3 = 1$,
- 896 • The parameters $\theta_1, \theta_2, \theta_3, \theta'_1, \theta'_2, \theta'_3$ are chosen so that all six experts are constant functions.
897 Let us define:

$$\begin{aligned} \mathcal{E}(\cdot; \theta_1) &= A_1, & \mathcal{E}(\cdot; \theta'_1) &= A_2, \\ \mathcal{E}(\cdot; \theta_2) &= B_1, & \mathcal{E}(\cdot; \theta'_2) &= B_2, \\ \mathcal{E}(\cdot; \theta_3) &= C_1, & \mathcal{E}(\cdot; \theta'_3) &= C_2. \end{aligned} \quad (79)$$

898 We now select the biases b_i, b'_i and constants A_j, B_j, C_j such that the models satisfy $\mathcal{D}(\cdot; \phi) =$
 899 $\mathcal{D}(\cdot; \phi')$, even though there exists no transformation as described in Theorem B.5 that maps ϕ to ϕ' .
 900 Our goal is to ensure that

$$\begin{aligned} & \frac{e^{-x+b_1}}{e^{-x+b_1} + e^{b_2} + e^{x+b_3}} \cdot A_1 + \frac{e^{b_2}}{e^{-x+b_1} + e^{b_2} + e^{x+b_3}} \cdot B_1 + \frac{e^{x+b_3}}{e^{-x+b_1} + e^{b_2} + e^{x+b_3}} \cdot C_1 \\ &= \frac{e^{-x+b'_1}}{e^{-x+b'_1} + e^{b'_2} + e^{x+b'_3}} \cdot A_2 + \frac{e^{b'_2}}{e^{-x+b'_1} + e^{b'_2} + e^{x+b'_3}} \cdot B_2 + \frac{e^{x+b'_3}}{e^{-x+b'_1} + e^{b'_2} + e^{x+b'_3}} \cdot C_2. \end{aligned} \quad (80)$$

901 For convenience, we introduce the following shorthand:

$$\begin{aligned} e^{b_1} &= X_1, & e^{b'_1} &= X_2, \\ e^{b_2} &= Y_1, & e^{b'_2} &= Y_2, \\ e^{b_3} &= Z_1, & e^{b'_3} &= Z_2. \end{aligned} \quad (81)$$

902 With this notation, Equation (80) becomes:

$$\begin{aligned} & \frac{e^{-x}X_1}{e^{-x}X_1 + Y_1 + e^xZ_1} \cdot A_1 + \frac{Y_1}{e^{-x}X_1 + Y_1 + e^xZ_1} \cdot B_1 + \frac{e^xZ_1}{e^{-x}X_1 + Y_1 + e^xZ_1} \cdot C_1 \\ &= \frac{e^{-x}X_2}{e^{-x}X_2 + Y_2 + e^xZ_2} \cdot A_2 + \frac{Y_2}{e^{-x}X_2 + Y_2 + e^xZ_2} \cdot B_2 + \frac{e^xZ_2}{e^{-x}X_2 + Y_2 + e^xZ_2} \cdot C_2, \end{aligned} \quad (82)$$

903 which is equivalent to:

$$\begin{aligned} & (e^{-x}X_1A_1 + Y_1B_1 + e^xZ_1C_1)(e^{-x}X_2 + Y_2 + e^xZ_2) \\ &= (e^{-x}X_2A_2 + Y_2B_2 + e^xZ_2C_2)(e^{-x}X_1 + Y_1 + e^xZ_1). \end{aligned} \quad (83)$$

904 By equating the coefficients of $e^{-2x}, e^{-x}, 1, e^x, e^{2x}$, we derive the following system:

$$\begin{aligned} e^{-2x} &: & X_1X_2A_1 &= & X_1X_2A_2, \\ e^{2x} &: & Z_1Z_2C_1 &= & Z_1Z_2C_2, \\ e^x &: & Y_1Z_2B_1 + Z_1Y_2C_1 &= & Y_1Z_2C_2 + Z_1Y_2B_2, \\ e^{-x} &: & Y_1X_2B_1 + X_1Y_2A_1 &= & Y_1X_2A_2 + X_1Y_2B_2, \\ 1 &: & X_1Z_2A_1 + Z_1X_2C_1 + Y_1Y_2B_1 &= & X_1Z_2C_2 + Z_1X_2A_2 + Y_1Y_2B_2. \end{aligned} \quad (84)$$

905 Setting $A_1 = A_2 = A$ and $C_1 = C_2 = C$ satisfies the equations involving e^{-2x} and e^{2x} automati-
 906 cally. Removing those, we simplify the system to:

$$\begin{aligned} e^x &: & Y_1Z_2B_1 + Z_1Y_2C &= & Y_1Z_2C + Z_1Y_2B_2, \\ e^{-x} &: & Y_1X_2B_1 + X_1Y_2A &= & Y_1X_2A + X_1Y_2B_2, \\ 1 &: & X_1Z_2A + Z_1X_2C + Y_1Y_2B_1 &= & X_1Z_2C + Z_1X_2A + Y_1Y_2B_2. \end{aligned} \quad (85)$$

907 Solving the equations for A and C , assuming $Z_1Y_2 \neq Y_1Z_2$ and $X_1Y_2 \neq Y_1X_2$, gives:

$$\begin{aligned} A &= \frac{X_1Y_2B_2 - Y_1X_2B_1}{X_1Y_2 - Y_1X_2}, \\ C &= \frac{Z_1Y_2B_2 - Y_1Z_2B_1}{Z_1Y_2 - Y_1Z_2}. \end{aligned} \quad (86)$$

908 Substituting into the constant term equation in (85) gives:

$$Y_1Y_2(B_1 - B_2) = (C - A)(X_1Z_2 - Z_1X_2). \quad (87)$$

909 Now computing $A - C$:

$$A - C = \frac{X_1Y_2B_2 - Y_1X_2B_1}{X_1Y_2 - Y_1X_2} - \frac{Z_1Y_2B_2 - Y_1Z_2B_1}{Z_1Y_2 - Y_1Z_2}$$

$$= \frac{Y_1 Y_2 (B_1 - B_2) (X_1 Z_2 - Z_1 X_2)}{(X_1 Y_2 - Y_1 X_2) (Z_1 Y_2 - Y_1 Z_2)}. \quad (88)$$

910 Plugging into Equation (87), we find:

$$Y_1 Y_2 (B_1 - B_2) = - \frac{Y_1 Y_2 (B_1 - B_2) (X_1 Z_2 - Z_1 X_2)}{(X_1 Y_2 - Y_1 X_2) (Z_1 Y_2 - Y_1 Z_2)} (X_1 Z_2 - Z_1 X_2). \quad (89)$$

911 Assuming $B_1 \neq B_2$ and $Y_1 Y_2 \neq 0$, we divide both sides to obtain:

$$(X_1 Y_2 - Y_1 X_2) (Y_1 Z_2 - Z_1 Y_2) = (X_1 Z_2 - Z_1 X_2)^2. \quad (90)$$

912 While this equation can be solved in general, it suffices to exhibit a concrete solution. Consider:

$$\begin{aligned} (X_1, X_2) &= (1, 2), \\ (Y_1, Y_2) &= (3, 5), \\ (Z_1, Z_2) &= (2, 3). \end{aligned} \quad (91)$$

913 With these values, B_1 and B_2 may be freely chosen. These assignments define parameter configura-
914 tions ϕ and ϕ' for which $\mathcal{D}(\cdot; \phi) = \mathcal{D}(\cdot; \phi')$, yet no transformation described in Theorem B.5 maps ϕ
915 to ϕ' .

916 C Results on Mixture-of-Experts with Dense Gating

917 C.1 Strongly Distinctness Property

918 The following definition formalizes the notion of *strong distinctness*, which will be used in the
919 statement of Theorem C.4.

920 **Definition C.1** (Strongly distinct). Two functions $f, g: X \rightarrow Y$ are said to be *strongly distinct* if the
921 set $\{x \in X : f(x) \neq g(x)\}$ is dense in X .

922 **Example C.2.** We present several examples to illustrate the concept of strong distinctness.

- 923 • Two distinct polynomials on \mathbb{R}^n or \mathbb{C}^n are strongly distinct.
- 924 • Two distinct holomorphic functions are strongly distinct.
- 925 • In contrast, two distinct locally affine functions are not necessarily strongly distinct. For
926 instance:

927 – Let $f_1, f_2: \mathbb{R} \rightarrow \mathbb{R}$ be defined as:

$$f_1(x) = \begin{cases} 0 & \text{if } x < 0, \\ x & \text{if } x \geq 0, \end{cases} \quad f_2(x) = 1. \quad (92)$$

928 Then f_1 and f_2 are strongly distinct.

929 – Let $g_1, g_2: \mathbb{R} \rightarrow \mathbb{R}$ be defined as:

$$g_1(x) = \begin{cases} 0 & \text{if } x < 0, \\ x & \text{if } x \geq 0, \end{cases} \quad g_2(x) = 0. \quad (93)$$

930 Here, g_1 and g_2 are distinct but not strongly distinct, since they coincide on $(-\infty, 0)$.

931 We now define a certain class of subsets of \mathbb{R}^d . Given

$$\{W_i, b_i\}_{i=1}^n \in (\mathbb{R}^d \times \mathbb{R})^n, \quad (94)$$

932 we define the set

$$\Omega(\{W_i, b_i\}_{i=1}^n) := \{x \in \mathbb{R}^d : \{W_i x + b_i\}_{i=1}^n \text{ are pairwise distinct}\}. \quad (95)$$

933 The following result provides a sufficient condition on the gating weights under which the Top- k
934 operator is capable of selecting any subset of k experts.

935 **Proposition C.3.** Suppose that the collection $\{W_i\}_{i=1}^n$ satisfies the condition that the set $\{W^{(G,i-1)} -$
 936 $W^{(G,i)}\}_{i=2}^n$ is linearly independent in \mathbb{R}^d . Then, for every subset $\mathcal{A} \subseteq \{1, \dots, n\}$ of size k , there
 937 exists a point $x \in \Omega(\{W_i, b_i\}_{i=1}^n)$ such that

$$\text{Top-}k((W_i x + b_i)_{i=1}^n) = \mathcal{A}. \quad (96)$$

938 *Proof.* Without loss of generality, assume $\mathcal{A} = \{1, \dots, k\}$. To prove the existence of $x \in \Omega$ such
 939 that

$$\text{Top-}k((W_i x + b_i)_{i=1}^n) = \{1, \dots, k\}, \quad (97)$$

940 it suffices to find $x \in \mathbb{R}^d$ such that

$$W_1 x + b_1 > W_2 x + b_2 > \dots > W_n x + b_n. \quad (98)$$

941 We can strengthen this to require:

$$\begin{aligned} (W_1 x + b_1) - (W_2 x + b_2) &= 1, \\ (W_2 x + b_2) - (W_3 x + b_3) &= 1, \\ &\vdots \\ (W_{n-1} x + b_{n-1}) - (W_n x + b_n) &= 1. \end{aligned} \quad (99)$$

942 This is equivalent to the following system of linear equations:

$$\begin{aligned} (W_1 - W_2)x &= 1 - (b_1 - b_2), \\ (W_2 - W_3)x &= 1 - (b_2 - b_3), \\ &\vdots \\ (W_{n-1} - W_n)x &= 1 - (b_{n-1} - b_n). \end{aligned} \quad (100)$$

943 Since the vectors $\{W_{i-1} - W_i\}_{i=2}^n$ are linearly independent by assumption, this system has a unique
 944 solution. Hence, there exists $x \in \mathbb{R}^d$ satisfying Equation (100). \square

945 Proposition C.3 will be invoked in the proof of Theorem C.4. A justification for the linear indepen-
 946 dence assumption is provided in Remark C.8.

947 C.2 Functional Equivalence in Mixture-of-Experts with Sparse Gating

948 We establish a functional equivalence theorem for the Sparse Mixture-of-Experts (SMoE) architecture,
 949 in parallel with the result for MoE given in Theorem B.5. However, our analysis is limited to the case
 950 $k > 1$, as the $k = 1$ setting leads to singularities that disrupt the general structure of equivalence. A
 951 detailed explanation for this exclusion is provided in Remark C.9.

952 **Theorem C.4** (Functional equivalence in Mixture-of-Experts with Sparse Gating). Suppose ϕ, ϕ'
 953 define the same SMoE function, i.e. $\mathcal{S}(\cdot; \phi) = \mathcal{S}(\cdot; \phi')$. Assume that the following conditions hold:

- 954 1. Both $\{\mathcal{E}(\cdot; \theta_i)\}_{i=1}^n$ and $\{\mathcal{E}(\cdot; \theta'_i)\}_{i=1}^{n'}$ consist of pairwise strongly distinct functions;
- 955 2. Both $\{W_{i-1} - W_i\}_{i=2}^n$ and $\{W'_{i-1} - W'_i\}_{i=2}^{n'}$ are linear independent subsets of \mathbb{R}^d .

956 Then $n = n'$, and there exists $g = (c_W, c_b, \tau) \in G(n)$ such that for all $i = 1, \dots, n$, we have
 957 $W'_i = W_{\tau(i)} + c_W$, $b'_i = b_{\tau(i)} + c_b$, and $\mathcal{E}(x; \theta'_i) = \mathcal{E}(x; \theta_{\tau(i)})$ for all $x \in \Omega(\{W_i, b_i\}_{i=1}^n)$ such that
 958 $\tau(i) \in \text{Top-}k((W_i x + b_i)_{i=1}^n)$.

959 Before presenting the proof of Theorem C.4, we begin with two remarks.

960 **Remark C.5.** Suppose $n = n'$ and there exist $\tau \in S_n$, $c_W \in \mathbb{R}^d$, and $c_b \in \mathbb{R}$ such that for every
 961 $i = 1, \dots, n$, we have

$$W'_i = W_{\tau(i)} + c_W, \quad b'_i = b_{\tau(i)} + c_b. \quad (101)$$

962 Then the two sets $\Omega(\{W_i, b_i\}_{i=1}^n)$ and $\Omega(\{W'_i, b'_i\}_{i=1}^n)$ coincide. Furthermore, for any x in this set,
 963 the condition $\tau(i) \in \text{Top-}k((W_i x + b_i)_{i=1}^n)$ holds if and only if $i \in \text{Top-}k((W'_i x + b'_i)_{i=1}^n)$.

964 **Remark C.6.** It is easy to verify that Assumption 2 in Theorem C.4 implies Assumption 2 in
 965 Theorem B.5.

966 *Proof.* To enhance clarity, we begin by outlining the main steps of the proof at a high level:

- 967 1. Formulate the identity $\mathcal{S}(\cdot; \phi) = \mathcal{S}(\cdot; \phi')$ explicitly and introduce simplified notation to
 968 streamline the presentation.
- 969 2. Partition the input space into regions where the Top- k operator selects the same index set,
 970 and within which all expert functions are affine.
- 971 3. Demonstrate the desired equivalence for a fixed subset of experts. The core technique is to
 972 invoke the MoE equivalence result given in Theorem B.5.
- 973 4. Generalize the argument to establish the equivalence across all experts.

974 We now proceed with the detailed derivation and proofs for each of the outlined steps.

975 **Step 1.** Given that $\mathcal{S}(\cdot; \phi) = \mathcal{S}(\cdot; \phi')$, it follows that

$$\begin{aligned} & \sum_{i \in \text{Top-}k((W_i x + b_i)_{i=1}^n)} \text{softmax}_i \left(\{W_i x + b_i\}_{i \in \text{Top-}k((W_i x + b_i)_{i=1}^n)} \right) \cdot \mathcal{E}(x; \theta_i) \\ &= \sum_{i \in \text{Top-}k((W'_i x + b'_i)_{i=1}^{n'})} \text{softmax}_i \left(\{W'_i x + b'_i\}_{i \in \text{Top-}k((W'_i x + b'_i)_{i=1}^{n'})} \right) \cdot \mathcal{E}(x; \theta'_i), \end{aligned} \quad (102)$$

976 for all $x \in \mathbb{R}^d$. To simplify notation, define

$$\mathcal{E}_i(\cdot) = \mathcal{E}(\cdot; \theta_i), \quad \text{and} \quad \mathcal{E}'_i(\cdot) = \mathcal{E}(\cdot; \theta'_i). \quad (103)$$

977 Using this notation, we can rewrite Equation (102) more compactly as:

$$\begin{aligned} & \sum_{i \in \text{Top-}k((W_i x + b_i)_{i=1}^n)} \text{softmax}_i \left(\{W_i x + b_i\}_{i \in \text{Top-}k((W_i x + b_i)_{i=1}^n)} \right) \cdot \mathcal{E}_i(x) \\ &= \sum_{i \in \text{Top-}k((W'_i x + b'_i)_{i=1}^{n'})} \text{softmax}_i \left(\{W'_i x + b'_i\}_{i \in \text{Top-}k((W'_i x + b'_i)_{i=1}^{n'})} \right) \cdot \mathcal{E}'_i(x). \end{aligned} \quad (104)$$

978 **Step 2.** We begin by highlighting two key observations:

- 979 • Assumption 2 guarantees that the parameter pairs $\{W_i, b_i\}$ are pairwise distinct for $i =$
 980 $1, \dots, n$, and likewise, $\{W'_i, b'_i\}$ are pairwise distinct for $i = 1, \dots, n'$. By Proposition A.1,
 981 the set

$$\Omega_1 = \Omega(\{W_i, b_i\}_{i=1}^n) \cap \Omega(\{W'_i, b'_i\}_{i=1}^{n'}), \quad (105)$$

982 is open and dense in \mathbb{R}^d . For any $x \in \Omega_1$, the values $\{W_i x + b_i\}$ and $\{W'_i x + b'_i\}$ are
 983 pairwise distinct. Moreover, for each $x \in \Omega_1$, there exists an open neighborhood around x
 984 in which both Top- k selections remain fixed.

- 985 • From the analysis in Appendix B.1, there exists a set $\Omega_2 \subset \mathbb{R}^d$, also open and dense, such
 986 that for every $x \in \Omega_2$, all expert functions \mathcal{E}_i and \mathcal{E}'_j are affine in a neighborhood of x .

987 Taking the intersection $\Omega = \Omega_1 \cap \Omega_2$, we obtain a set that is still open and dense. Within Ω ,
 988 both the Top- k selections and the expert functions remain locally constant and affine, respectively.
 989 Consequently, there exists a collection of open sets $\{U_i\}_{i \in I}$ that cover Ω , such that

$$\Omega = \bigcup_{i \in I} U_i, \quad (106)$$

990 and on each U_i , the expert functions \mathcal{E}_i and \mathcal{E}'_j are affine, and the Top- k selections do not vary.

991 **Step 3.** Let U be an arbitrary open set from the cover described in Equation (106). Without loss of
 992 generality, we may relabel the indices such that both Top- k maps are constantly equal to $\{1, \dots, k\}$
 993 throughout U . Under this reindexing, Equation (104) simplifies to

$$\begin{aligned} \sum_{i=1}^k \text{softmax}_i \left(\{W_i x + b_i\}_{i=1}^k \right) \cdot \mathcal{E}_i(x) \\ = \sum_{i=1}^k \text{softmax}_i \left(\{W'_i x + b'_i\}_{i=1}^k \right) \cdot \mathcal{E}'_i(x) \quad \text{for all } x \in U. \end{aligned} \quad (107)$$

994 According to Assumption 1, the expert functions \mathcal{E}_i are strongly distinct, which implies that they
 995 remain distinct on the open set U . The same conclusion applies to the functions \mathcal{E}'_i . Therefore, the
 996 assumptions of Theorem B.5 are satisfied on U , and Equation (107) falls within its scope. As a
 997 consequence, there exist constants $c_W \in \mathbb{R}^d$ and $c_b \in \mathbb{R}$ such that, up to reindexing, we have for all
 998 $i = 1, \dots, k$,

$$W'_i = W_i + c_W, \quad b'_i = b_i + c_b, \quad (108)$$

999 and furthermore, $\mathcal{E}_i = \mathcal{E}'_i$ on U .

1000 **Step 4.** For any index $m \in \{3, 4, \dots, n\}$, we invoke Proposition C.3 to select an open set V_1 from
 1001 the cover in Equation (106) such that both indices 1 and k appear in $T(V_1)$. Restricting Equation (104)
 1002 to V_1 and applying Theorem C.4, we conclude that there exist indices $1 \leq t_1, s_1 \leq n'$ such that

$$W_1 - W_m = W'_{t_1} - W'_{s_1}. \quad (109)$$

1003 Applying the same reasoning with indices 2 and m , we obtain $1 \leq t_2, s_2 \leq n'$ such that

$$W_2 - W_m = W'_{t_2} - W'_{s_2}. \quad (110)$$

1004 Subtracting Equation (110) from Equation (109) yields

$$W'_1 - W'_2 = W_1 - W_2 = (W_1 - W_m) - (W_2 - W_m) = (W'_{t_1} - W'_{s_1}) - (W'_{t_2} - W'_{s_2}). \quad (111)$$

1005 Due to the linear independence guaranteed by Assumption 2, this identity can only hold if $t_1 = 1$,
 1006 $t_2 = 2$, and $s_1 = s_2$. Denoting this shared index by $\tau(m)$, i.e., $\tau(m) = s_1 = s_2$, we then have

$$W_1 - W_m = W'_1 - W'_{\tau(m)}, \quad (112)$$

1007 which implies

$$W'_{\tau(m)} - W_m = W'_1 - W_1 = c_W. \quad (113)$$

1008 Similarly, we deduce

$$b'_{\tau(m)} - b_m = b'_1 - b_1 = c_b. \quad (114)$$

1009 Since m ranges over $\{3, 4, \dots, n\}$, the corresponding values of $\tau(m)$ must be distinct. If not, suppose
 1010 that $\tau(m) = \tau(m')$ for some $m \neq m'$, which would imply

$$W_m - W_{m'} = W'_{\tau(m)} - W'_{\tau(m')} = 0, \quad (115)$$

1011 contradicting Assumption 3.

1012 Applying a symmetric argument to the parameters of $\mathcal{S}(\cdot; \phi')$, we conclude that $n = n'$. Therefore,
 1013 there exists a permutation τ of $\{1, \dots, n\}$ such that

$$W'_i = W_{\tau(i)} + c_W, \quad b'_i = b_{\tau(i)} + c_b. \quad (116)$$

1014 Finally, the analysis above shows that for any $x \in \Omega(\{W_i, b_i\}_{i=1}^n)$ with $\tau(i) \in \text{Top-}k((W_i x +$
 1015 $b_i)_{i=1}^n)$ —i.e., when index i is selected by the Top- k operator in \mathcal{S} —we have

$$\mathcal{E}_{\tau(i)}(x) = \mathcal{E}'_i(x). \quad (117)$$

1016 This completes the proof of Theorem C.4. \square

Remark C.7. While Theorem C.4 is conceptually analogous to Theorem B.5, it is crucial to recognize that establishing the result for SMoE involves substantially greater technical complexity. The main challenge arises from the Top- k operator, which introduces discontinuities by dynamically changing the subset of active experts in a manner that depends intricately on the input. This input-dependent behavior complicates the analysis and makes the theoretical treatment significantly more delicate.

Remark C.8 (Rationale behind the assumptions in Theorem C.4). We begin by restating the two assumptions made in Theorem C.4:

1. Both $\{\mathcal{E}(\cdot; \theta_i)\}_{i=1}^n$ and $\{\mathcal{E}'(\cdot; \theta'_i)\}_{i=1}^{n'}$ consist of pairwise strongly distinct functions;
2. Both $\{W_{i-1} - W_i\}_{i=2}^n$ and $\{W'_{i-1} - W'_i\}_{i=2}^{n'}$ are linear independent subsets of \mathbb{R}^d .

These assumptions are strictly stronger than those required in Theorem B.5. We now discuss their necessity and implications in greater detail.

Assumption 1. This condition arises primarily from the behavior of the Top- k operator, which induces input-dependent expert selection. As a result, the functional contribution of an expert is restricted to regions where it is actively selected by the gating mechanism. Outside these regions, the expert can behave arbitrarily without influencing the output. Therefore, if the experts are merely pairwise distinct—rather than pairwise strongly distinct—it becomes possible for different sets of expert functions to yield identical overall behavior when restricted to their respective activation domains. This ambiguity underscores the necessity of strong distinctness to ensure functional identifiability in the SMoE setting.

Assumption 2. In practice, the number of experts n is typically much smaller than the input (token) dimension D . As a result, the collections $\{W^{(G,i-1)} - W_i\}_{i=2}^n$ and $\{W'^{(G,i-1)} - W'_i\}_{i=2}^{n'}$ are generically linearly independent. However, when linear dependence arises, it can prevent certain expert pairs from ever being simultaneously selected by the gating mechanism across all possible inputs. This limitation introduces singular symmetries: different parameter configurations that yield functionally identical outputs but cannot be related via the equivalence structure defined in Theorem C.4.

To illustrate this phenomenon concretely, consider an example with $n = 4$ and $k = 2$, and let $\mathcal{E}_1, \mathcal{E}_2, \mathcal{E}_3, \mathcal{E}_4$ denote arbitrary expert functions. Define two SMoE models \mathcal{S}_1 and \mathcal{S}_2 , whose gating logits are $(-2x, -x, x, 2x)$ and $(-3x, -2x, 2x, 3x)$, respectively. The resulting functions take the form:

$$\mathcal{S}_1(x) = \begin{cases} \text{softmax}_1(-2x, -x) \cdot \mathcal{E}_1(x) + \text{softmax}_2(-2x, -x) \cdot \mathcal{E}_2(x) & \text{if } x < 0, \\ \text{softmax}_1(x, 2x) \cdot \mathcal{E}_3(x) + \text{softmax}_2(x, 2x) \cdot \mathcal{E}_4(x) & \text{if } x > 0, \end{cases} \quad (118)$$

and

$$\mathcal{S}_2(x) = \begin{cases} \text{softmax}_1(-3x, -2x) \cdot \mathcal{E}_1(x) + \text{softmax}_2(-3x, -2x) \cdot \mathcal{E}_2(x) & \text{if } x < 0, \\ \text{softmax}_1(2x, 3x) \cdot \mathcal{E}_3(x) + \text{softmax}_2(2x, 3x) \cdot \mathcal{E}_4(x) & \text{if } x > 0. \end{cases} \quad (119)$$

It is easy to verify that $\mathcal{S}_1(x) = \mathcal{S}_2(x)$ for all $x \in \mathbb{R} \setminus \{0\}$, where the gating logits are pairwise distinct and the Top- k selections remain constant. However, no transformation of the form specified in Theorem C.4 maps one configuration to the other. This demonstrates how singular symmetries can arise in the SMoE architecture, even when functional outputs coincide on an open dense subset of the input domain.

Remark C.9 (The case of $k = 1$). In the special case where $k = 1$, the SMoE function in Equation (21) simplifies to

$$\mathcal{S}(x; \{W_i, b_i, \theta_i\}_{i=1}^n) = \mathcal{E}(x; \theta_i), \quad (120)$$

where the index i is determined by

$$i = \underset{i=1, \dots, n}{\operatorname{argmax}} (W_i x + b_i). \quad (121)$$

In this setting, the Top-1 gating mechanism selects only the expert with the highest score, and the softmax reduces to a one-hot distribution with a single nonzero entry equal to 1.

Beyond the standard $G(n)$ symmetry acting on the expert parameters, the SMoE architecture with $k = 1$ exhibits an additional, nontrivial invariance under the action of the multiplicative group $\mathbb{R}_{>0}$. Specifically, for any scalar $c > 0$, we have

$$\mathcal{S}(x; \{W_i, b_i, \theta_i\}_{i=1}^n) = \mathcal{S}(x; \{cW_i, cb_i, \theta_i\}_{i=1}^n), \quad (122)$$

as the argmax used to select the active expert remains invariant under uniform positive scaling:

$$\operatorname{argmax}_{i=1, \dots, n} (W_i x + b_i) = \operatorname{argmax}_{i=1, \dots, n} (cW_i x + cb_i), \quad (123)$$

for all $x \in \Omega(\{W_i, b_i\}_{i=1}^n)$. Furthermore, since only one expert is active at any given input, the architecture does not involve explicit interactions between experts. This leads to a more complex symmetry structure, including hidden and continuous transformations, which complicates theoretical analysis. For this reason, we exclude the case $k = 1$ from our main results and leave its investigation to future work.

D Technical Details for Sections 5 and 6

D.1 Proof for the Sufficiency of Permutation Invariance in LMC of MoE

The following result establishes that translation transformations do not affect the barrier loss between two parameter configurations.

Proposition D.1. *Let $h = (c_W, c_b, \text{id}_n) \in G(n)$, where $\text{id}_n \in S_n$ denotes the identity permutation. For any $\phi_A, \phi_B \in \Phi(n)$, the barrier loss remains invariant under translation:*

$$B(\phi_A, \phi_B) = B(\phi_A, h\phi_B). \quad (124)$$

That is, the translation components do not contribute to the barrier loss as defined in Equation (1).

Proof. Recall that

$$B(\phi_A, \phi_B) = \sup_{t \in [0, 1]} [\mathcal{L}(t\phi_A + (1-t)\phi_B) - (t\mathcal{L}(\phi_A) + (1-t)\mathcal{L}(\phi_B))]. \quad (125)$$

For $t \in [0, 1]$, denote $h_t = (tc_W, tc_b, \text{id}_n) \in G(n)$. For $\phi_A, \phi_B \in \Phi(n)$ such that

$$\phi_A = (W_i^A, b_i^A, \theta_i^A)_{i=1, \dots, n} \quad (126)$$

$$\phi_B = (W_i^B, b_i^B, \theta_i^B)_{i=1, \dots, n}, \quad (127)$$

we have:

$$\begin{aligned} & t\phi_A + (1-t)(h\phi_B) \\ &= t(W_i^A, b_i^A, \theta_i^A)_{i=1, \dots, n} + (1-t)(W_i^B + c_W, b_i^B + c_b, \theta_i^B)_{i=1, \dots, n} \\ &= (tW_i^A + (1-t)W_i^B + (1-t)c_W, tb_i^A + (1-t)b_i^B + (1-t)c_b, t\theta_i^A + (1-t)\theta_i^B)_{i=1, \dots, n} \\ &= (h_{1-t})(t\phi_A + (1-t)\phi_B). \end{aligned} \quad (128)$$

Since the group action of $G(n)$ on $\Phi(n)$ preserves the functionality of the MoE function \mathcal{D} , we have

$$\mathcal{D}(\cdot; \phi_B) = \mathcal{D}(\cdot; h\phi_B), \quad (129)$$

and

$$\begin{aligned} & \mathcal{D}(\cdot; t\phi_A + (1-t)(h\phi_B)) \\ &= \mathcal{D}(\cdot; (h_{1-t})(t\phi_A + (1-t)\phi_B)) \\ &= \mathcal{D}(\cdot; t\phi_A + (1-t)\phi_B). \end{aligned} \quad (130)$$

These observations lead to

$$\mathcal{L}(\phi_B) = \mathcal{L}(h\phi_B), \quad (131)$$

$$\mathcal{L}(t\phi_A + (1-t)(h\phi_B)) = \mathcal{L}(t\phi_A + (1-t)\phi_B). \quad (132)$$

1080 Thus, for $t \in [0, 1]$, we have

$$\begin{aligned}
& B(\phi_A, h\phi_B) \\
&= \sup_{t \in [0, 1]} \left[\mathcal{L}(t\phi_A + (1-t)(h\phi_B)) - (t\mathcal{L}(\phi_A) + (1-t)\mathcal{L}(h\phi_B)) \right] \\
&= \sup_{t \in [0, 1]} \left[\mathcal{L}(t\phi_A + (1-t)\phi_B) - (t\mathcal{L}(\phi_A) + (1-t)\mathcal{L}(\phi_B)) \right] \\
&= B(\phi_A, \phi_B).
\end{aligned} \tag{133}$$

1081 Therefore, the proof is finished. \square

1082 D.2 Proof of the Permutation-Invariant Property for Equation (11)

1083 **Proposition D.2.** *The cost function defined in Method 2 of Section 5.2 is permutation-invariant with*
1084 *respect to the hidden units of the experts in the MoE models.*

1085 *Proof.* For each expert i in the first Mixture-of-Experts (MoE) model ϕ , let $A_i \in \mathbb{R}^{h \times d}$ and $u_i \in \mathbb{R}^h$
1086 be the weight matrix and bias for the first layer, and $B_i \in \mathbb{R}^{d \times h}$ and $v_i \in \mathbb{R}^d$ be the weight
1087 matrix and bias for the second layer. Define the augmented matrices $\tilde{A}_i = [A_i, u_i] \in \mathbb{R}^{h \times (d+1)}$
1088 and $\tilde{B}_i = [B_i, v_i] \in \mathbb{R}^{d \times (h+1)}$. Similarly, for each expert j in the second MoE model ϕ' , define
1089 A'_j, u'_j, B'_j, v'_j , and the augmented matrices $\tilde{A}'_j, \tilde{B}'_j$ accordingly. We define the Gram matrices
1090 $\tilde{A}_i^T \tilde{A}_i \in \mathbb{R}^{(d+1) \times (d+1)}$ and $\tilde{B}_i \tilde{B}_i^T \in \mathbb{R}^{d \times d}$ for each expert i in ϕ , and similarly for ϕ' .

1091 To establish permutation invariance, consider a permutation matrix $P \in \mathbb{R}^{h \times h}$ applied to the hidden
1092 units of expert i in ϕ . For the first layer, the permuted augmented weight matrix is $P\tilde{A}_i$. For the
1093 second layer, the permutation affects the columns of B_i within $\tilde{B}_i = [B_i, v_i]$, so the transformed
1094 matrix is $\tilde{B}_i \begin{bmatrix} P^T & 0 \\ 0 & 1 \end{bmatrix}$, leaving the bias v_i unchanged. We now verify the invariance of the Gram
1095 matrices under this permutation. For the first layer:

$$(P\tilde{A}_i)^T (P\tilde{A}_i) = \tilde{A}_i^T P^T P \tilde{A}_i = \tilde{A}_i^T I \tilde{A}_i = \tilde{A}_i^T \tilde{A}_i, \tag{134}$$

1096 since $P^T P = I$, the $h \times h$ identity matrix. For the second layer:

$$\begin{aligned}
& \left(\tilde{B}_i \begin{bmatrix} P^T & 0 \\ 0 & 1 \end{bmatrix} \right) \left(\tilde{B}_i \begin{bmatrix} P^T & 0 \\ 0 & 1 \end{bmatrix} \right)^T = \tilde{B}_i \begin{bmatrix} P^T & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} P & 0 \\ 0 & 1 \end{bmatrix} \tilde{B}_i^T \\
&= \tilde{B}_i \begin{bmatrix} P^T P & 0 \\ 0 & 1 \end{bmatrix} \tilde{B}_i^T \\
&= \tilde{B}_i \begin{bmatrix} I & 0 \\ 0 & 1 \end{bmatrix} \tilde{B}_i^T \\
&= \tilde{B}_i \tilde{B}_i^T,
\end{aligned} \tag{135}$$

1097 where $\begin{bmatrix} I & 0 \\ 0 & 1 \end{bmatrix}$ is the $(h+1) \times (h+1)$ identity matrix. This confirms that both $\tilde{A}_i^T \tilde{A}_i$ and $\tilde{B}_i \tilde{B}_i^T$ are
1098 invariant under permutations of the hidden units.

1099 The cost matrix $C \in \mathbb{R}^{n \times n}$ for the Linear Assignment Problem (LAP) is $C = \{C_{i,j}\}_{1 \leq i \leq n, 1 \leq j \leq n}$,
1100 where

$$C_{i,j} = \left(\left\| \tilde{A}_i^T \tilde{A}_i - (\tilde{A}'_j)^T \tilde{A}'_j \right\|_F^2 + \left\| \tilde{B}_i \tilde{B}_i^T - \tilde{B}'_j (\tilde{B}'_j)^T \right\|_F^2 \right)^{\frac{1}{2}} \tag{136}$$

1101 where $\|\cdot\|_F$ denotes the Frobenius norm. Since the Gram matrices $(\tilde{A}_i)^T \tilde{A}_i$ and $\tilde{B}_i (\tilde{B}_i)^T$ (and
1102 their counterparts in ϕ') are permutation-invariant, their differences and the Frobenius norms of
1103 these differences are also invariant. Consequently, $C_{i,j}$ remains unchanged under permutations of the
1104 hidden units in either ϕ or ϕ' .

1105 Thus, the cost function is permutation-invariant with respect to the hidden units of the experts,
1106 completing the proof. \square

1107 D.3 Weight Matching Algorithm for Mixture-of-Experts

Algorithm 1 Weight Matching for Mixture-of-Experts

Input: MoE model weights $\phi = (W_i, b_i, \theta_i)_{i=1, \dots, n}$, $\phi' = (W'_i, b'_i, \theta'_i)_{i=1, \dots, n}$
Output: Permutation τ for experts, and permutations $\{P_i\}_{i=1}^n$ for hidden units
 % Step 1: Match experts' order using two methods
for method in {gate, expert} **do**
 Compute cost matrix C
 Solve LAP to obtain expert permutation τ_{method}
end for
 % The two candidate expert orderings τ_{gate} and τ_{expert} are obtained
 % Step 2: Align internal weights of matched expert pairs
for method in {gate, expert} **do**
 for $i = 1$ to n **do**
 Compute P_i by applying Weight Matching to θ_i and $\theta'_{\tau_{\text{method}}(i)}$
 end for
end for
return $(\tau_{\text{gate}}, (\{P_i\}_{i=1}^n)_{\text{gate}}), (\tau_{\text{expert}}, (\{P_i\}_{i=1}^n)_{\text{expert}})$

1108 D.4 Formal formulation of DeepSeekMoE

1109 For completeness, we present the notation and formulations of three variants of MoE models: Dense
 1110 MoE, Sparse MoE (SMoE), and DeepSeekMoE.

1111 **Expert.** Given d denoting the input dimension. We define an *Expert* as a function $\mathcal{E}(\cdot; \theta): \mathbb{R}^d \rightarrow \mathbb{R}^d$,
 1112 parameterized by $\theta \in \mathbb{R}^e$ with $e \in \mathbb{N}$ is the total number of trainable parameters of \mathcal{E} . In this work,
 1113 we consider each expert $\mathcal{E}(\cdot; \theta)$ to be a feedforward neural network with ReLU activation functions.
 1114 Unless stated otherwise, all experts are assumed to share the same architecture.

1115 **Mixture of Expert with dense gating.** Given n denoting the number of experts, we define *Mixture-*
 1116 *of-Experts with dense gating* as a function $\mathcal{D}: \mathbb{R}^d \rightarrow \mathbb{R}^d$ such that

$$\mathcal{D}(x; \{W_i, b_i, \theta_i\}_{i=1}^n) = \sum_{i=1}^n \text{softmax}_i(s_1(x), \dots, s_n(x)) \mathcal{E}(x; \theta_i), \quad (137)$$

1117 where $s_i(x) = W_i x + b_i$ determines the contribution of each expert to the final output. Here, θ_i
 1118 are the parameters of the i^{th} expert. The function $s = (s_1, \dots, s_n)$ is called the *gating score*, and is
 1119 parameterized as $s(\cdot; \{W_i, b_i\}_{i=1}^n)$ with $(W_i, b_i) \in \mathbb{R}^d \times \mathbb{R}$ are the corresponding *gating* parameters.

1120 **Mixture-of-Experts with sparse gating.** Given a positive integer $k \leq n$ denoting the number of
 1121 activated experts, define the Top- k map by $\text{Top-}k(z) = \{i_1, \dots, i_k\}$ for $z = (z_1, \dots, z_n) \in \mathbb{R}^n$,
 1122 where i_1, \dots, i_k are the indices corresponding to the k largest components of x . In the event of
 1123 ties, we select smaller indices first. We define a *Mixture-of-Experts with sparse gating* (SMoE) as a
 1124 function $\mathcal{S}: \mathbb{R}^d \rightarrow \mathbb{R}^d$ as follows. For $x \in \mathbb{R}^d$, let $T(x) = \text{Top-}k(s(x))$, then define:

$$\mathcal{S}(x; \{W_i, b_i, \theta_i\}_{i=1}^n) = \sum_{i \in T(x)} \text{softmax}_i((s_i(x))_{i \in T(x)}) \cdot \mathcal{E}(x; \theta_i) \quad (138)$$

1125 In other words, the Top- k selects the k highest-scoring experts used to compute the output.

1126 **DeepseekMoE with shared and routed experts.** The DeepseekMoE architecture extends the sparse
 1127 Mixture-of-Experts by incorporating both shared experts that are always active and routed experts
 1128 that are sparsely activated based on the input. Given positive integers n_s and n_r denoting the number
 1129 of shared and routed experts, respectively ($n = n_s + n_r$ as the total number of experts), and a positive
 1130 integer $k_r \leq n_r$ denoting the number of activated routed experts, the DeepseekMoE layer is defined
 1131 as a function $\mathcal{M}: \mathbb{R}^d \rightarrow \mathbb{R}^d$ such that

$$\mathcal{M}(x; \{\theta_i^{(s)}\}_{i=1}^{n_s}, \{W_i, b_i, \theta_i\}_{i=1}^{n_r}) = \sum_{i=1}^{n_s} \mathcal{E}_i^{(s)}(x; \theta_i^{(s)}) + \sum_{i=1}^{n_r} g_i(x) \mathcal{E}(x; \theta_i^{(r)}), \quad (139)$$

1132 where:

- $\mathcal{E}_i^{(s)}(\cdot; \theta_i^{(s)}): \mathbb{R}^d \rightarrow \mathbb{R}^d$ are the shared experts, parameterized by $\theta_i^{(s)}$,
- $\mathcal{E}_i^{(r)}(\cdot; \theta_i^{(r)}): \mathbb{R}^d \rightarrow \mathbb{R}^d$ are the routed experts, parameterized by $\theta_i^{(r)}$,
- The gating values $g_i(x)$ are computed as follows:

$$\begin{aligned}\tilde{s}(x) &= \text{softmax}(s_1(x), \dots, s_{n_r}(x)) \in \mathbb{R}^{n_r}, \\ g_i(x) &= \begin{cases} \tilde{s}_i(x), & \text{if } i \in \text{Top-}k_r(\tilde{s}(x)) \\ 0, & \text{otherwise} \end{cases}\end{aligned}$$

This formulation allows the model to leverage a combination of always-active shared experts and input-dependent routed experts, enhancing efficiency and performance in large-scale neural networks.

E Impact of Feedforward Reinitialization on Pretrained Transformer Performance

We empirically investigate the sensitivity of pretrained Transformer models to targeted parameter reinitialization within their Feedforward Network (FFN) modules. The FFN, a critical component for non-linear transformations and feature representation within each Transformer block, was selectively reinitialized layer by layer to assess its functional contribution and the robustness of the overall model architecture. This analysis aims to identify the extent to which performance is dependent on the learned parameters of individual FFNs and to potentially highlight layers that are more critical to maintaining task proficiency. To this end, experiments were conducted on two representative model-task pairs:

1. Image Classification using a Vision Transformer (ViT-Base, patch size 16, image resolution 224) pretrained on ImageNet.
2. Language Modeling using a GPT-2 model (12-layer, hidden size 768) pretrained on Wiki-Text103.

For each model, the parameters of the FFN subcomponent within a single Transformer block were reinitialized using a standard initialization scheme (e.g., variance scaling), while all other model parameters, including attention weights and the FFNs in other layers, were kept frozen from the pretrained state. Following reinitialization, the models were evaluated on their respective tasks without any subsequent fine-tuning. This approach isolates the immediate effect of disrupting a single FFN on pretrained performance.

Figures 3 and 4 illustrate the performance degradation observed when reinitializing the FFN at different layers for the ViT and GPT-2 models, respectively. A striking observation is the significantly larger performance drop (increase in loss, decrease in accuracy) when reinitializing the FFN in the first layer (Layer 0) compared to any subsequent layer. While reinitializing FFNs in layers beyond the first does result in performance degradation, this impact is markedly less severe than at Layer 0. The results further reveal a non-uniform sensitivity profile among these subsequent layers, with FFN reinitialization in middle layers tending to induce a more pronounced performance decrement than in later layers. This pattern suggests a unique functional importance or sensitivity associated with the initial layer’s FFN, while the network’s architecture, particularly the presence of skip connections, appears to enable a greater degree of resilience to disruption in FFNs at deeper layers.

A key observation from these experiments, particularly for deeper and more complex models with skip connections, is that reinitializing a single FFN (except potentially the very first one) often does not lead to a catastrophic performance collapse. This phenomenon can be partially attributed to the presence of skip connections, which facilitate the unimpeded flow of information across layers, allowing features learned by other parts of the network to bypass the disrupted FFN. Furthermore, standard weight initialization schemes typically employ relatively small scales centered around zero, meaning the reinitialized FFN parameters introduce a perturbation that is initially small and balanced compared to the potentially large and specialized weights learned during pretraining. This combination of architectural properties (skip connections) and initialization characteristics helps the model maintain functionality, suggesting that the reinitialized state may remain within or close to the original optimization basin found during pretraining. Previous work [53] has highlighted how skip

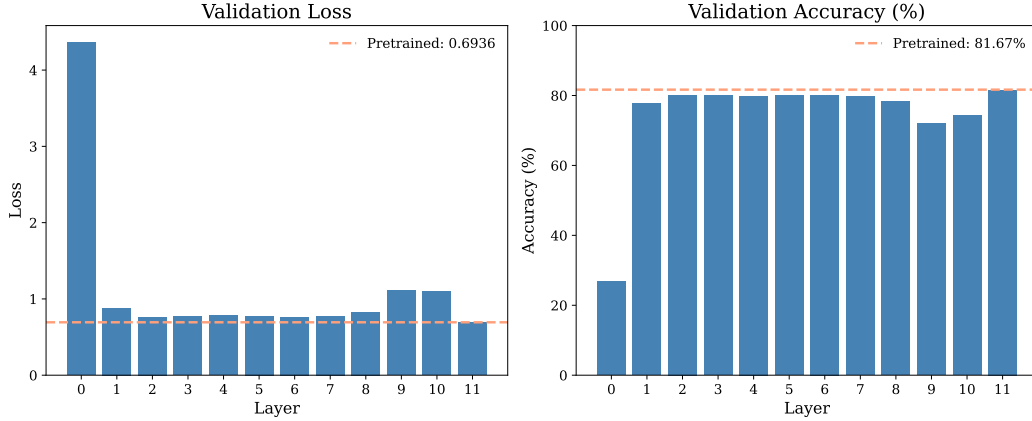


Figure 3: Performance degradation in ViT-Base on ImageNet due to FFN reinitialization at different layers.

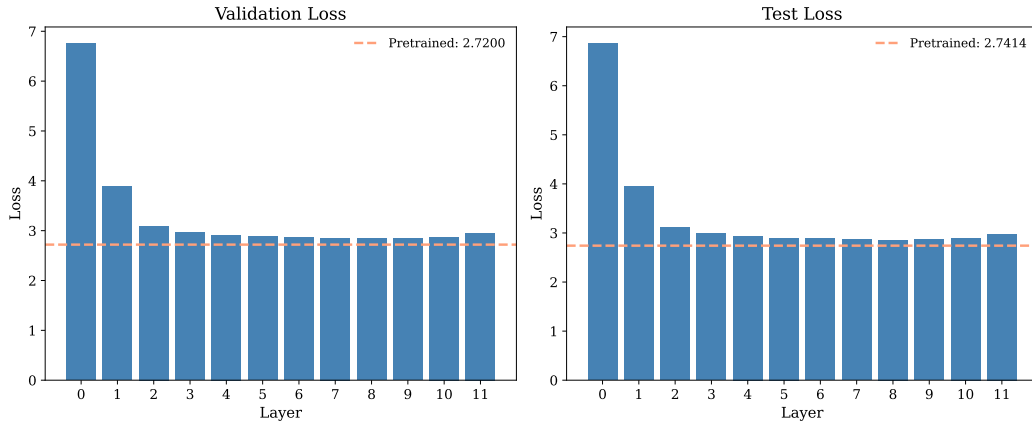


Figure 4: Effect of FFN reinitialization on GPT-2 perplexity across layers on WikiText103.

connections contribute to flatter minima and improve training stability in deep networks, which aligns with the observation that reinitialization causes less disruption than might intuitively be expected.

To further investigate the generality of this observation, we also conducted FFN reinitialization experiments on tasks commonly used for evaluating model robustness and optimization landscape properties, including text classification tasks (IMDB Review [57], AGNEWS [86], DBPedia [8]) and additional language modeling tasks (EnWik8 [38], Penn Treebank (Word Level) [58]). Table 4 presents a summary of the performance changes on these datasets after reinitializing an arbitrary single FFN layer compared to the original pretrained performance, and also shows performance after subsequent fine-tuning.

As shown in Table 4, while reinitializing a single FFN does lead to performance degradation on these tasks (e.g., loss increases by approximately 0.1 to 0.3), the magnitude of this change is relatively modest compared to the significant loss increases observed in experiments involving more widespread reinitialization or on larger tasks like GPT-2 on WikiText103 or ViT on ImageNet (where full model reinitialization can increase loss by several units). The relatively small performance drop suggests that for these specific datasets and this focused reinitialization strategy, the model’s state remains close to its original pretrained optimum. Consequently, metrics used to assess mode connectivity between the original and reinitialized states, such as the Naive Loss Barrier (shown in Table 4), yield very low values. This indicates that a simple linear path in parameter space between the original and reinitialized single-FFN models exhibits low loss values. However, due to the minor performance impact of the reinitialization itself, this low barrier value may primarily reflect the fact that the reinitialized state is already highly functional and located within the same or a very accessible

Table 4: Effect of single FFN reinitialization on various tasks. "Pretrained" indicates performance of the original model. "Reinitialize" shows performance after reinitializing a single FFN layer. "Finetune" shows performance after fine-tuning the reinitialized model. "Naive Loss Barrier" refers to the loss barrier observed between two fine-tuned models initialized with different random seeds and trained with different batch orders.

Dataset	Pretrained Loss	Pretrained Accuracy	Reinitialize Loss	Reinitialize Accuracy	Finetune Loss	Finetune Accuracy	Naive Loss Barrier
IMDB Review	0.3724	87.50	0.6452	68.75	0.4284	87.50	0.0019
AGNEWS	0.2849	90.50	0.6006	76.64	0.2986	90.73	0.0007
DBPedia	0.1896	93.75	0.3017	87.50	0.2036	94.47	0.0003
EnWik8	0.9684	-	1.3053	-	0.9506	-	0.0082
Penn Treebank	4.5246	-	6.9699	-	4.4639	-	-0.0454

optimization basin, rather than necessarily implying a broadly flat landscape between drastically different high-performing modes. Therefore, while consistent with observations of flat minima facilitated by skip connections [53], these specific reinitialization experiments on these datasets may not provide deep insights into complex mode connectivity far from the original optimum.

In summary, these findings highlight the heterogeneous functional roles of FFNs across different layers in Transformer models and demonstrate varying degrees of sensitivity to parameter reinitialization. The relative robustness observed for many layers, especially in deep models, is consistent with the structural benefits of skip connections and the properties of standard initialization. This layer-wise sensitivity analysis contributes to a better understanding of Transformer architecture and has potential implications for model compression techniques (e.g., identifying less critical FFNs), efficient fine-tuning strategies, and the design of conditional computation mechanisms like Mixture-of-Experts (MoE) routing.

F Experimental Details and Hyperparameters

We conduct a comprehensive evaluation of Linear Mode Connectivity (LMC) across a broad suite of vision and language modeling benchmarks. For vision tasks, our study includes MNIST, CIFAR-10, CIFAR-100, ImageNet-1k, as well as transfer learning scenarios from ImageNet-21k to CIFAR-10 and CIFAR-100. For language modeling, we utilize WikiText103 and the One Billion Word dataset. All experiments are performed using pretrained Transformer-based architectures, wherein all original model parameters are frozen and only the inserted Mixture-of-Experts (MoE) layers are subject to fine-tuning.

For vision tasks, we adopt Vision Transformer (ViT) backbones, while GPT-2 serves as the backbone for language modeling. Model configurations—including architecture depth, width, and tokenization context—are adjusted per dataset to reflect task-specific complexity and scale. Each expert module is architecturally aligned with the original feedforward network (MLP) layers of the pretrained model. To promote stable convergence, learning rates are independently tuned for each setting within the range of 10^{-4} to 10^{-1} .

MNIST We utilize the original MNIST grayscale images, each resized into a grid of non-overlapping patches with a patch size of 7. A lightweight Vision Transformer (ViT) model is employed, configured with an embedding dimension of 32 and a depth of 1 to 2 Transformer layers, depending on the specific setup. Each self-attention layer uses 4 attention heads to process the patch embeddings. Dropout is applied at a rate of 0.0 and the activation function throughout the network is `gelu`. The model is optimized using stochastic gradient descent (SGD) during both the pretraining and fine-tuning stages.

CIFAR-10. We utilize the original CIFAR-10 images with a patch size of 4. The ViT model is employed, with an embedding dimension of 128, 8 self-attention heads and 2 to 6 transformer layers. Dropout is applied at a rate of 0.0, and the `gelu` activation function is used. The SGD optimizer is employed during both pretraining and fine-tuning.

CIFAR-100. We resize CIFAR-100 images to 224x224, with a patch size of 16. The ViT model is employed, with an embedding dimension of 384, 6 self-attention heads and 6 to 12 transformer layers.

Dropout is applied at a rate of 0.0, and the `gelu` activation function is used. The SGD optimizer is employed during fine-tuning.

Imagenet21k→CIFAR10. We adopt the ViT-Small-Patch16-224 model [75], pretrained on ImageNet-21k and subsequently fine-tuned on CIFAR-10. The model consists of 12 layers, a hidden size of 384, an MLP size of 1536, and 6 attention heads, resulting in approximately 22.2M parameters. It employs a patch size and stride of 16. Dropout is disabled (set to 0.0), and the activation function is `gelu`. Stochastic Gradient Descent (SGD) is employed during fine-tuning.

Imagenet21k→CIFAR100. We adopt the ViT-Small-Patch16-224 model [75], pretrained on ImageNet-21k and subsequently fine-tuned on CIFAR-100. The model consists of 12 layers, a hidden size of 384, an MLP size of 1536, and 6 attention heads, resulting in approximately 22.2M parameters. It employs a patch size and stride of 16. The attention probability dropout and hidden layer dropout are set to 0.0, and the activation function is `gelu`. The SGD optimizer is employed during fine-tuning.

ImageNet-1k. We use the ViT-Base-Patch16-224 model [20], consisting of 12 Transformer layers, each with a hidden size of 768, MLP hidden dimension of 3072, and 12 self-attention heads. The total parameter count is approximately 86M. The encoder uses a patch size and stride of 16, and the `gelu` activation function is applied throughout. Both the attention dropout and hidden dropout rates are set to 0.0. Optimization is performed using the Adam optimizer during both pretraining and MoE-layer fine-tuning.

WikiText103. We adopt a GPT-2 model architecture consisting of 12 layers with 12 attention heads and an embedding dimension of 768. The context length and maximum position embeddings are set to 1024. The dropout rate is set to 0.1, and the activation function is `gelu`. The vocabulary size is 50,257. The model uses the standard GPT-2 initialization and layer normalization settings. The Adam optimizer is employed during both pretraining and fine-tuning.

One Billion Word. Similar to WikiText103, the GPT-2 model consists of 12 layers, 12 attention heads, and an embedding dimension of 768, but with a reduced context and position length of 256. The dropout rate and activation function remain consistent with the WikiText103 setup. The vocabulary size is extended to 793,470 to accommodate the dataset’s linguistic diversity. The Adam optimizer is employed during both pretraining and fine-tuning.

All experiments are executed on a single NVIDIA H100 GPU with 80GB of memory, except for the One Billion Word task, which utilizes two H100 GPUs. Due to the use of the JAX framework, approximately 75% of GPU memory (around 60GB) is pre-allocated by default. The number of CPU workers used in data loading and preprocessing is kept less than or equal to 10 across all experiments. For smaller-scale tasks such as MNIST, CIFAR-10, CIFAR-100, and transfer learning from ImageNet-21k, each experiment completes in under 30 minutes. For larger-scale tasks, the fine-tuning durations are approximately 15 hours for ImageNet-1k, 3.5 hours for WikiText103, and 4 hours for One Billion Word.

G Experimental Results

G.1 Verification of Linear Mode Connectivity across diverse configurations

To rigorously evaluate the robustness and generality of Linear Mode Connectivity (LMC) in expert-based Transformer architectures, we conducted systematic experiments across a wide range of Mixture-of-Experts (MoE) configurations. Our study includes three representative variants—vanilla MoE, SMoE, and DeepSeekMoE—which differ in both architectural structure and expert routing mechanisms. In all cases, we substituted the feed-forward network (FFN) of the *first* Transformer layer with an MoE module, while leaving the rest of the model architecture unchanged. All models were fine-tuned from identical random initializations to ensure consistency in comparison.

Table 5 summarizes the experimental conditions, covering diverse model depths (2, 6, 12 layers), task domains (including vision benchmarks such as MNIST, CIFAR-10, ImageNet21k→CIFAR-100, and language modeling with WikiText103), and expert configurations. SMoE employs a sparsity level of $k = 2$ active experts per token, whereas DeepSeekMoE further incorporates a shared expert ($s = 1$) to encourage parameter reuse and stability. Across nearly all configurations, LMC consistently emerges: linear interpolation between independently fine-tuned models (initialized identically) yields

Table 5: Experimental configurations for LMC analysis across MoE, SMoE, and DeepSeekMoE variants, evaluated on multiple datasets and settings. Datasets of the form $A \rightarrow B$ denote pretraining on A and finetuning on B . SMoE uses $k = 2$ active experts, while DeepSeekMoE uses $k = 2$ with an additional shared expert ($s = 1$). In all cases, the MLP in the *first* Transformer layer is replaced with an MoE module.

Method	Dataset	No. layers	No. experts	Figure
MoE	MNIST	1	[2, 4]	[5, 6]
		2	[2, 4]	[7, 8]
	CIFAR-10	2	[2, 4, 6]	[9, 10, 11]
		6	[2, 4, 6]	[12, 13, 14]
	CIFAR-100	6	[2, 4, 6]	[15, 16, 17]
	ImageNet-21k→CIFAR-10	12	[2, 4, 6]	[18, 19, 20]
	ImageNet-21k→CIFAR-100	12	[2, 4, 6]	[21, 22, 23]
	ImageNet-1k	12	[2, 4, 6, 8]	[24, 25, 26, 27]
	WikiText103	12	[2, 4, 6, 8]	[28, 29, 30, 31]
	One Billion Word	12	[2, 4, 6, 8]	[32, 33, 34, 35]
SMoE ($k = 2$)	MNIST	1	[4]	[36]
		2	[4]	[37]
	CIFAR-10	2	[4, 8]	[38, 39]
		6	[4, 8]	[40, 41]
	CIFAR-100	6	[4, 8]	[42, 43]
	ImageNet-21k→CIFAR-10	12	[4, 8]	[44, 45]
	ImageNet-21k→CIFAR-100	12	[4, 8]	[46, 47]
	ImageNet-1k	12	[4, 8, 16]	[48, 49, 50]
	WikiText103	12	[4, 8, 16]	[51, 52, 53]
	One Billion Word	12	[4, 8, 16]	[54, 55, 56]
DeepSeekMoE ($k = 2, s = 1$)	MNIST	1	[4]	[57]
		2	[4]	[58]
	CIFAR-10	2	[4, 8]	[59, 60]
		6	[4, 8]	[61, 62]
	CIFAR-100	6	[4, 8]	[63, 64]
	ImageNet-21k→CIFAR-10	12	[4, 8]	[65, 66]
	ImageNet-21k→CIFAR-100	12	[4, 8]	[67, 68]
	ImageNet-1k	12	[4, 8, 16]	[69, 70, 71]
	WikiText103	12	[4, 8, 16]	[72, 73, 74]
	One Billion Word	12	[4, 8, 16]	[75, 76, 77]

smooth loss trajectories, with no significant barriers. This suggests that expert-based models, even with dynamic routing and varying capacities, tend to converge to connected optima under matched training setups. The corresponding loss curves, linked in Table 5, reinforce this observation.

These results provide strong empirical evidence that LMC is a robust and general phenomenon in expert architectures, supporting the hypothesis that expert specialization and routing do not disrupt the connected geometry of the optimization landscape when alignment in initialization and architecture is preserved.

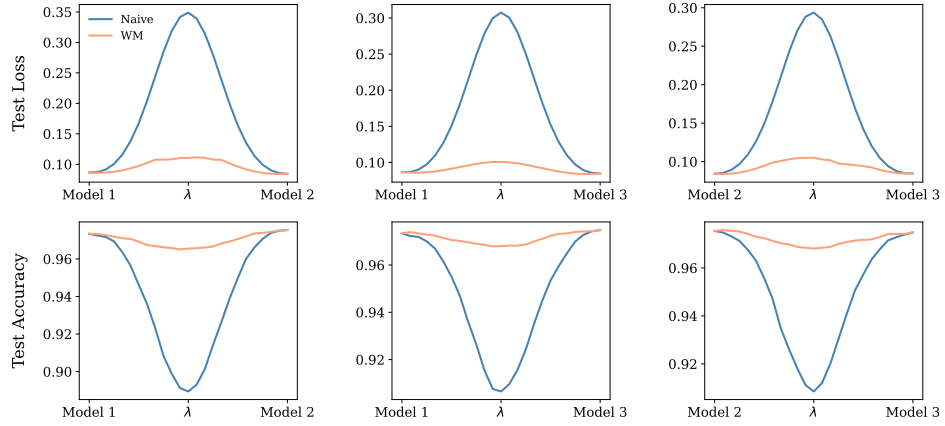


Figure 5: Linear Mode Connectivity for ViT-MoE on MNIST with 1 layer and 2 experts

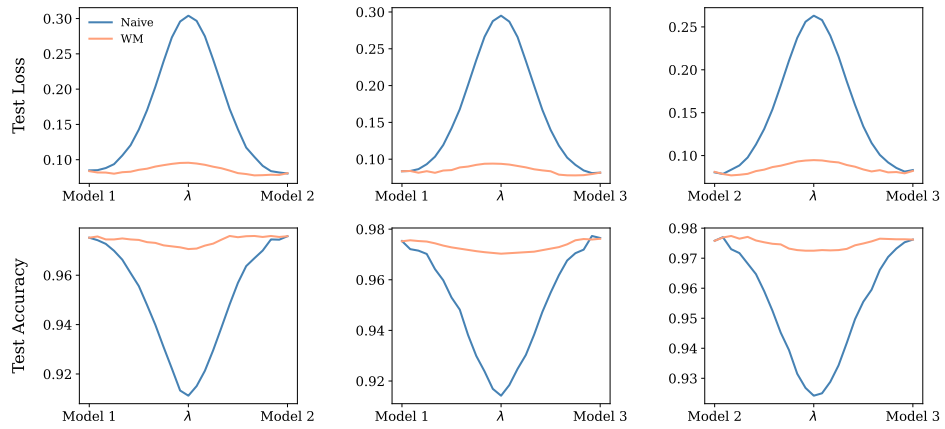


Figure 6: Linear Mode Connectivity for ViT-MoE on MNIST with 1 layer and 4 experts

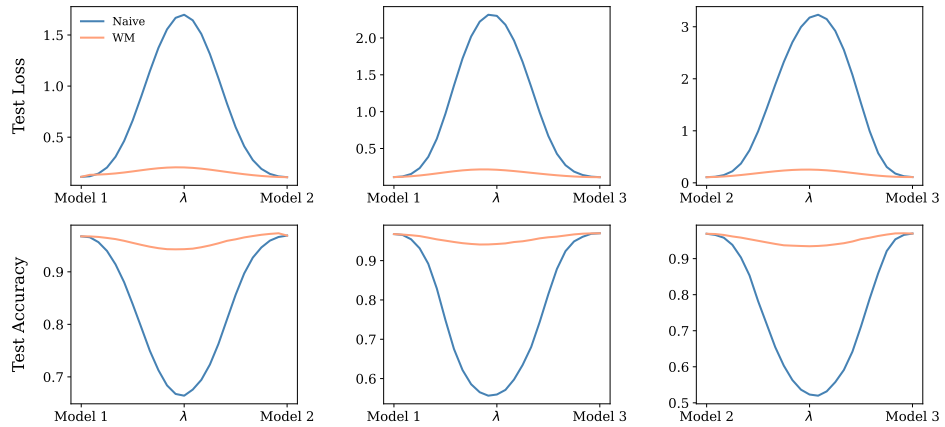


Figure 7: Linear Mode Connectivity for ViT-MoE on MNIST with 2 layers and 2 experts

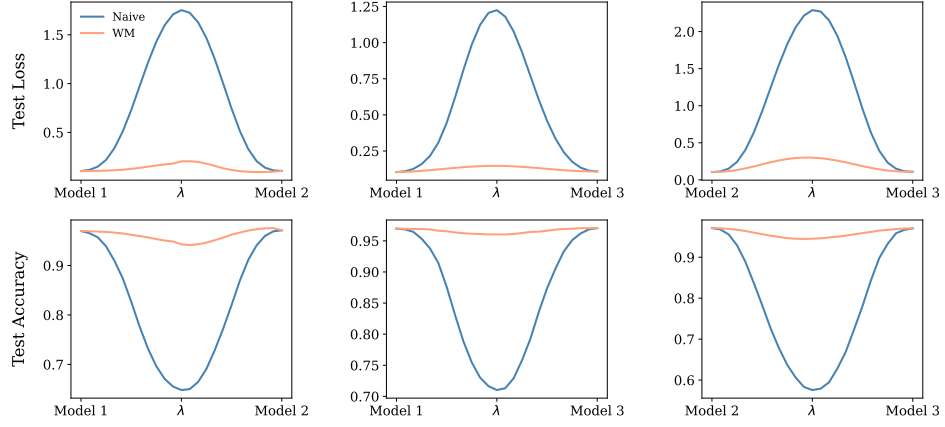


Figure 8: Linear Mode Connectivity for ViT-MoE on MNIST with 2 layers and 4 experts

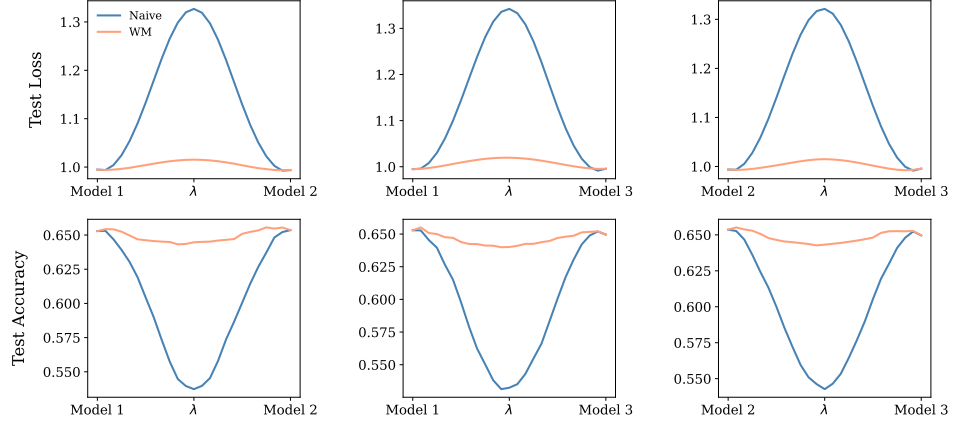


Figure 9: Linear Mode Connectivity for ViT-MoE on CIFAR-10 with 2 layers and 2 experts

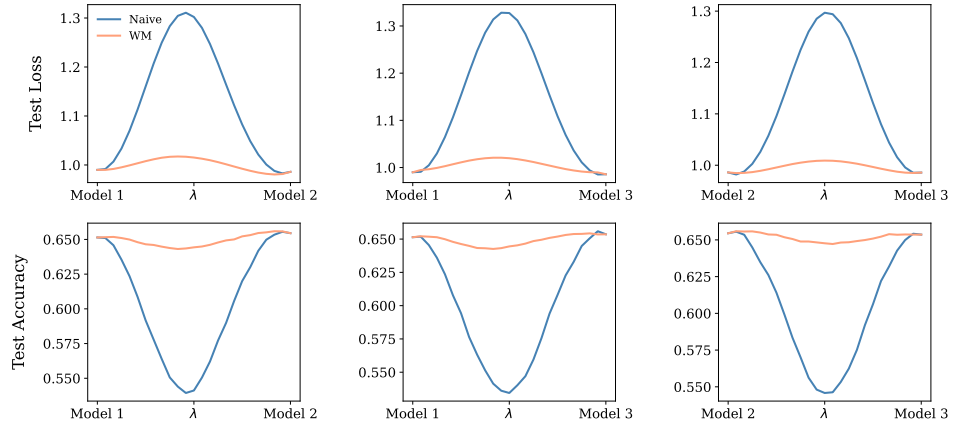


Figure 10: Linear Mode Connectivity for ViT-MoE on CIFAR-10 with 2 layers and 4 experts

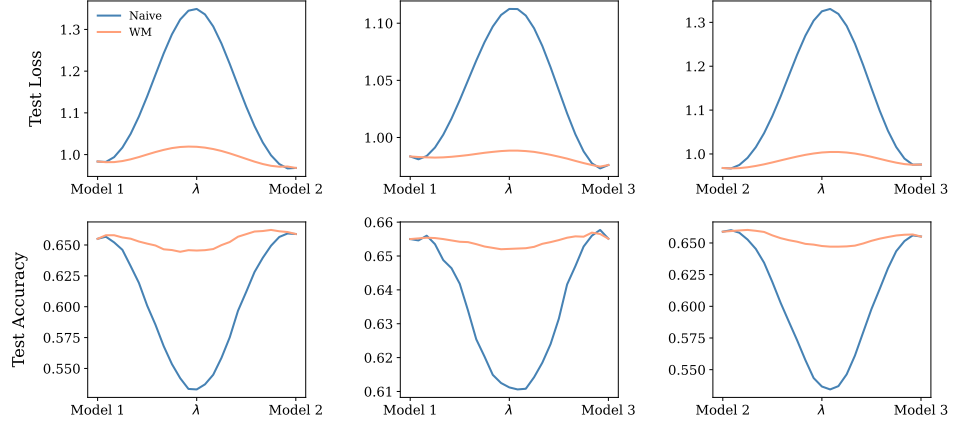


Figure 11: Linear Mode Connectivity for ViT-MoE on CIFAR-10 with 2 layers and 6 experts

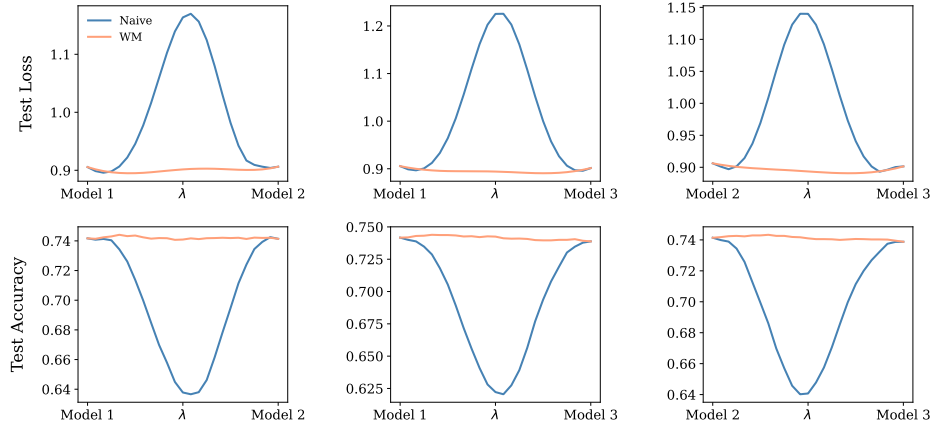


Figure 12: Linear Mode Connectivity for ViT-MoE on CIFAR-10 with 6 layers and 2 experts

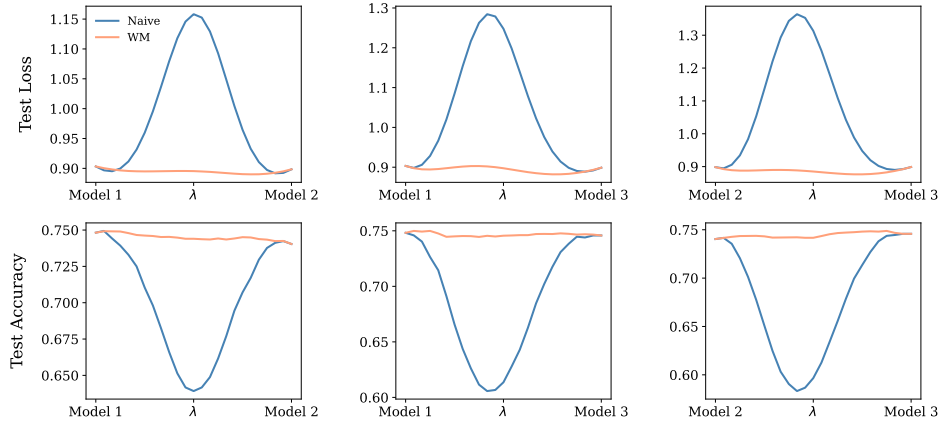


Figure 13: Linear Mode Connectivity for ViT-MoE on CIFAR-10 with 6 layers and 4 experts

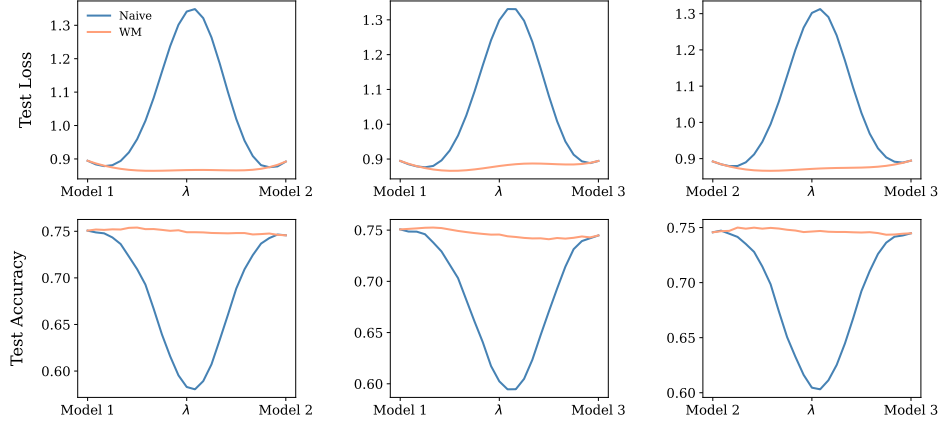


Figure 14: Linear Mode Connectivity for ViT-MoE on CIFAR-10 with 6 layers and 6 experts

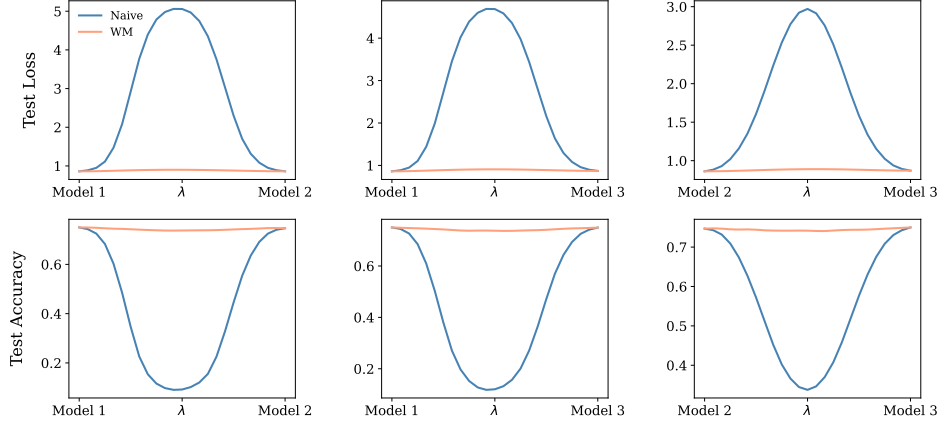


Figure 15: Linear Mode Connectivity for ViT-MoE on CIFAR-100 with 6 layers and 2 experts

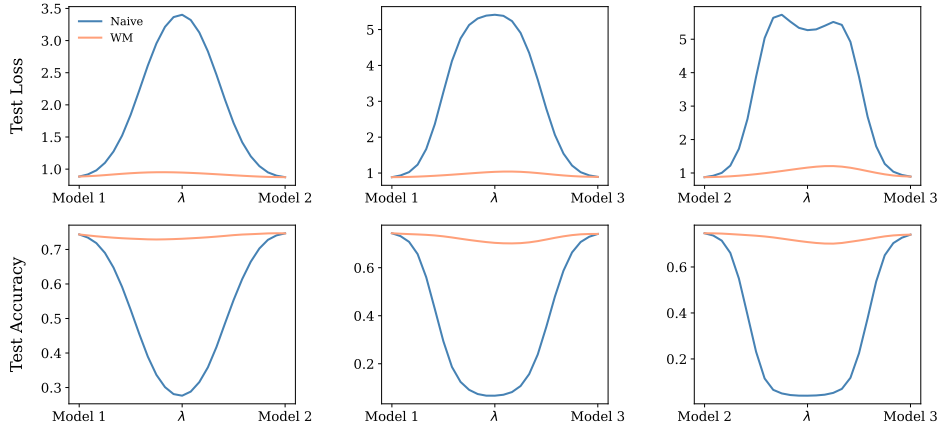


Figure 16: Linear Mode Connectivity for ViT-MoE on CIFAR-100 with 6 layers and 4 experts

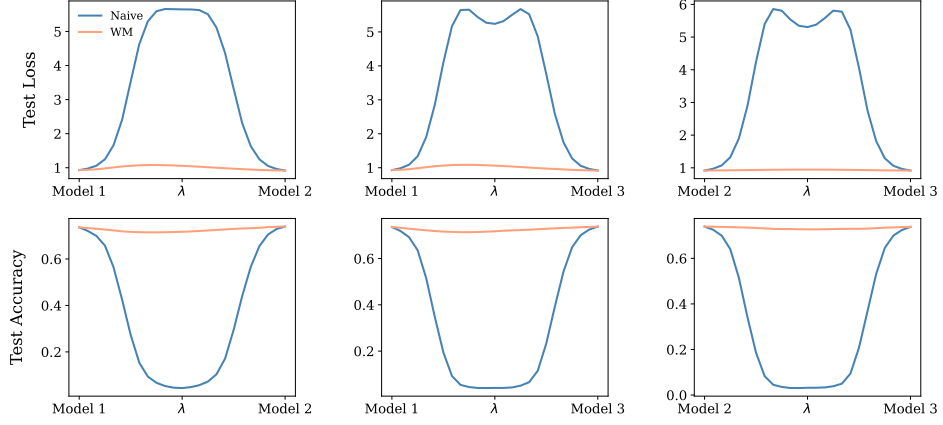


Figure 17: Linear Mode Connectivity for ViT-Moe on CIFAR-100 with 6 layers and 6 experts

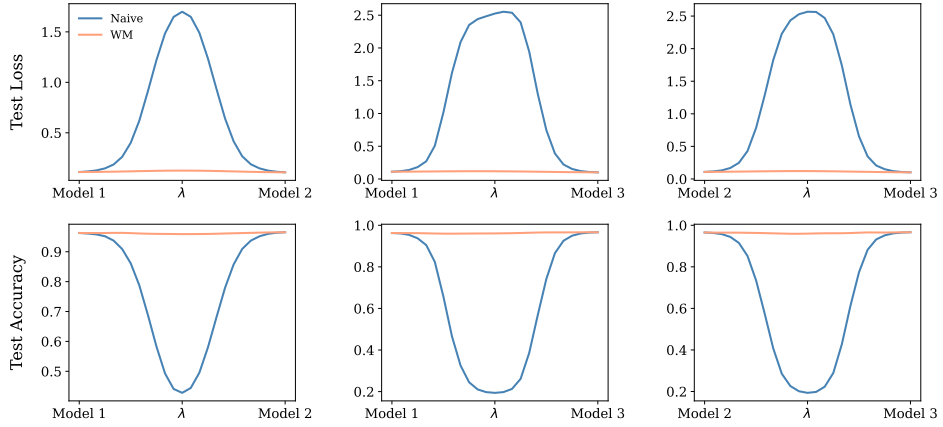


Figure 18: Linear Mode Connectivity for ViT-MoE on ImageNet-21k→CIFAR-10 with 12 layers and 2 experts

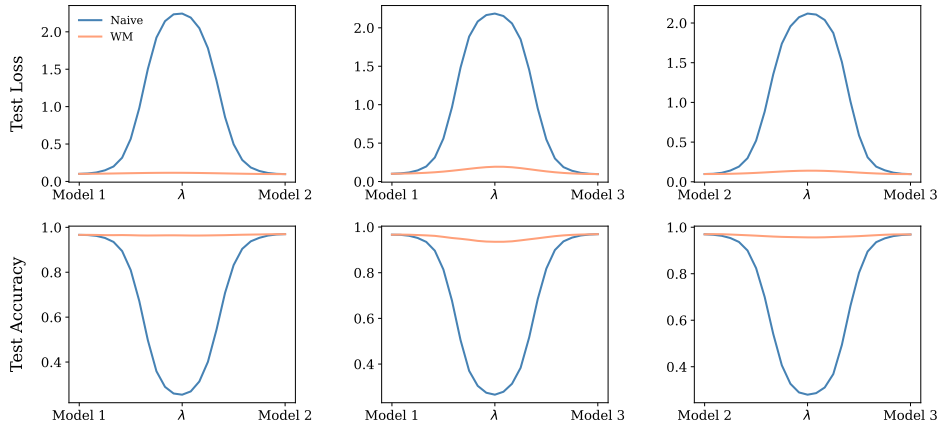


Figure 19: Linear Mode Connectivity for ViT-MoE on ImageNet-21k→CIFAR-10 with 12 layers and 4 experts

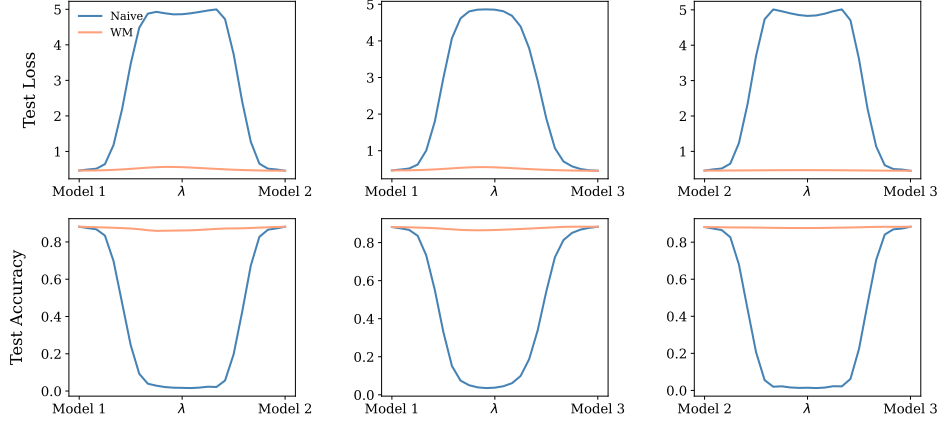


Figure 22: Linear Mode Connectivity for ViT-MoE on ImageNet-21k→CIFAR-100 with 12 layers and 4 experts

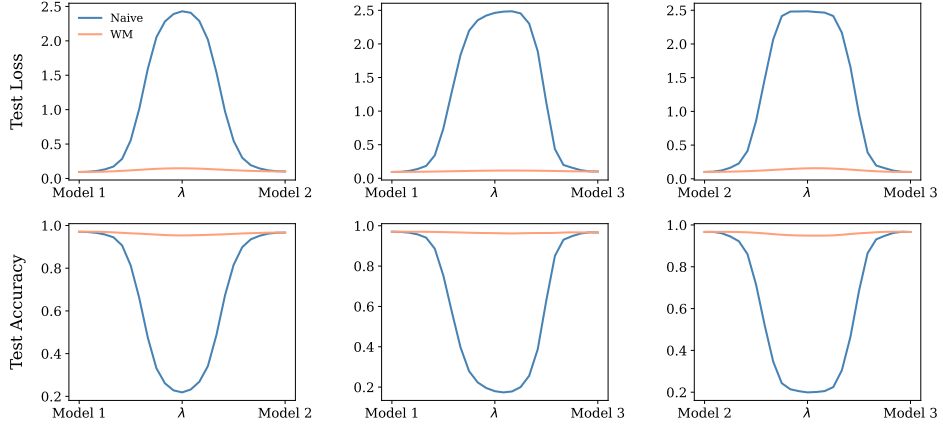


Figure 20: Linear Mode Connectivity for ViT-MoE on ImageNet-21k→CIFAR-100 with 12 layers and 6 experts

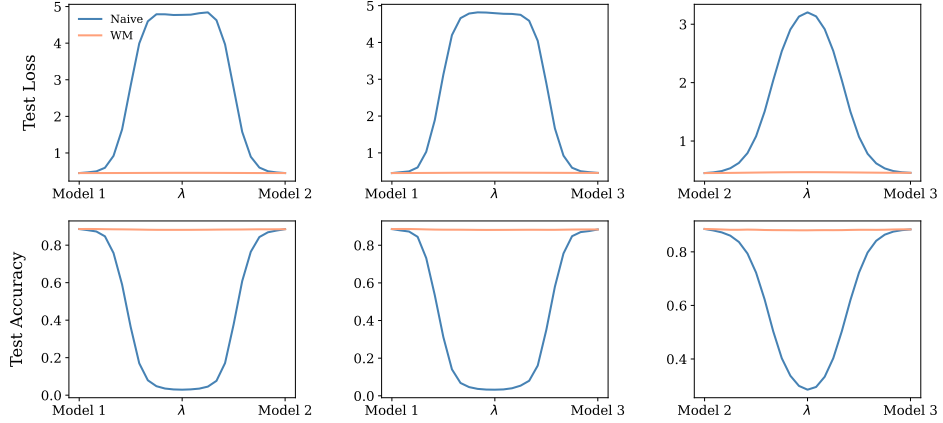


Figure 21: Linear Mode Connectivity for ViT-MoE on ImageNet-21k→CIFAR-100 with 12 layers and 2 experts

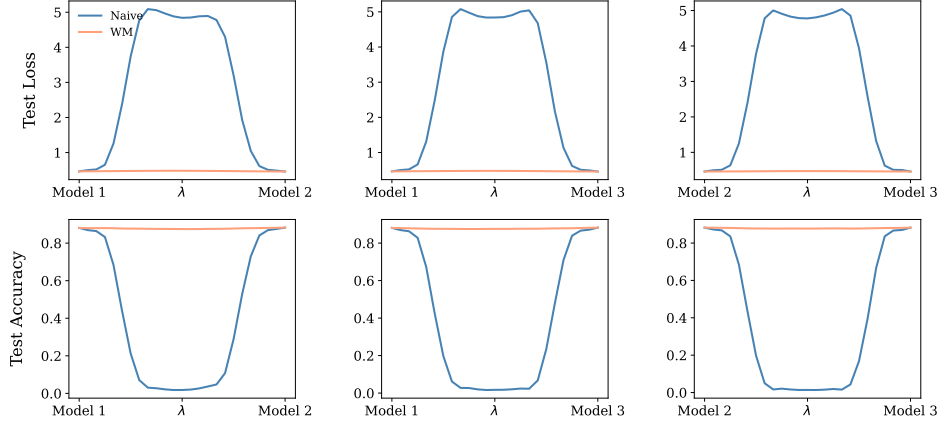


Figure 23: Linear Mode Connectivity for ViT-MoE with ImageNet-21k→CIFAR-100 with 12 layers and 6 experts

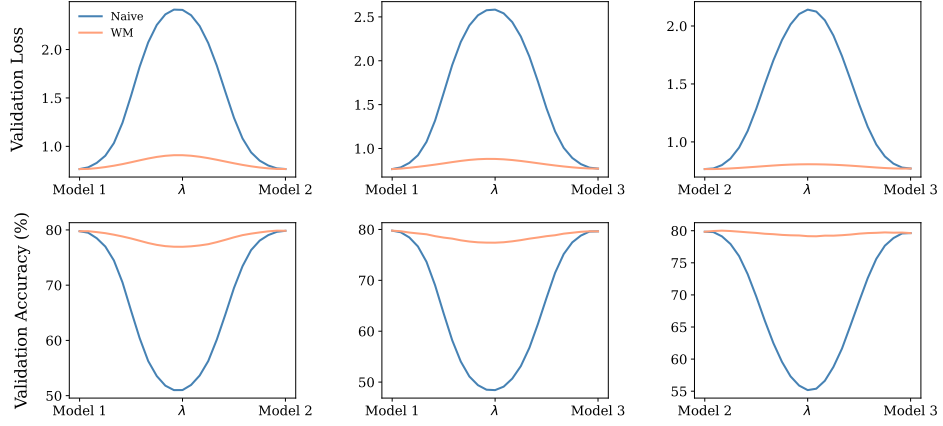


Figure 24: Linear Mode Connectivity for ViT-MoE on ImageNet-1k with 12 layers and 2 experts

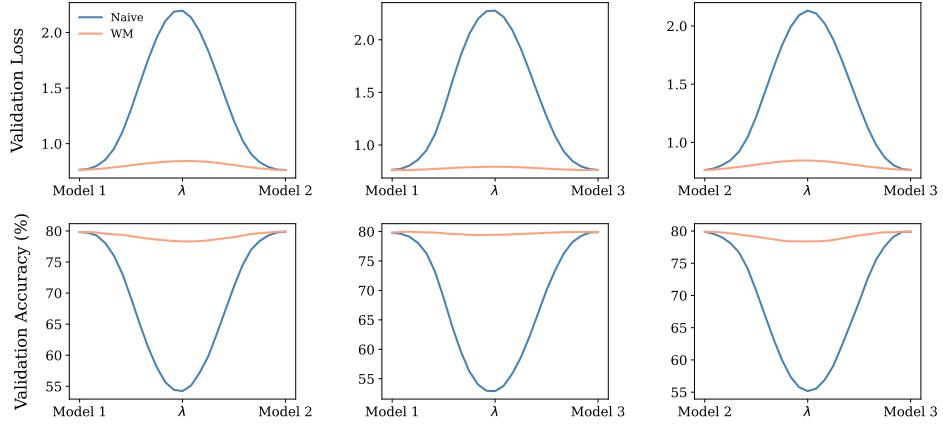


Figure 25: Linear Mode Connectivity for ViT-MoE on ImageNet-1k with 12 layers and 4 experts

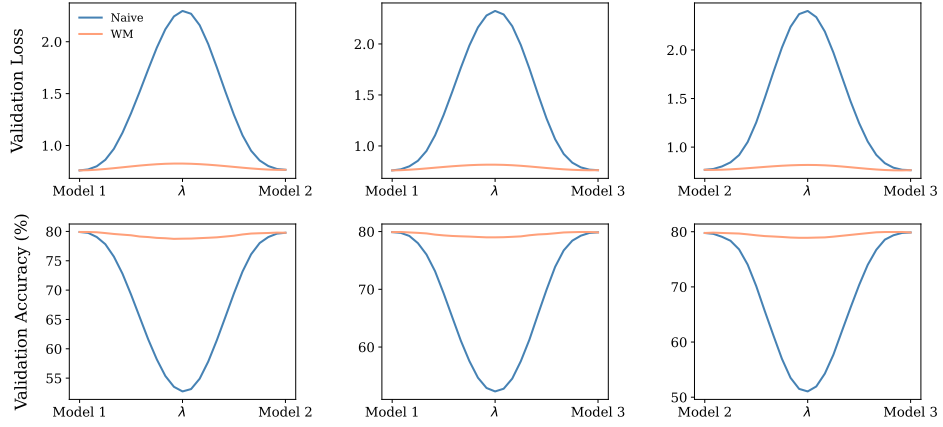


Figure 26: Linear Mode Connectivity for ViT-MoE on ImageNet-1k with 12 layers and 6 experts

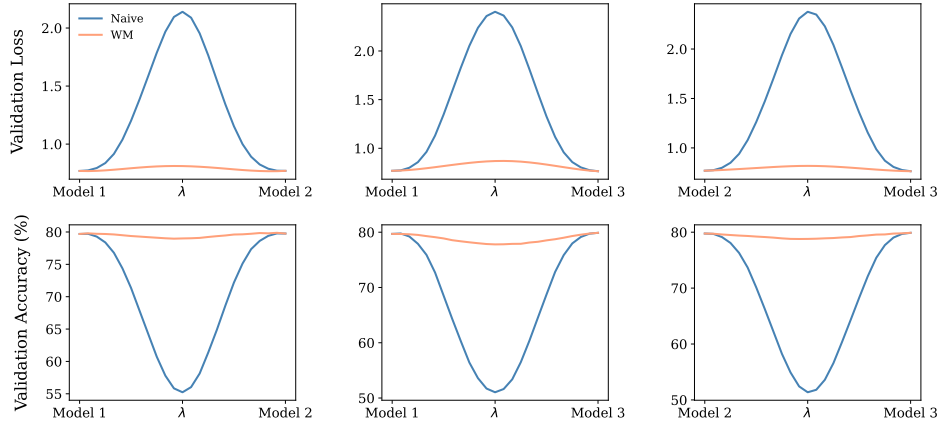


Figure 27: Linear Mode Connectivity for ViT-MoE on ImageNet-1k with 12 layers and 8 experts

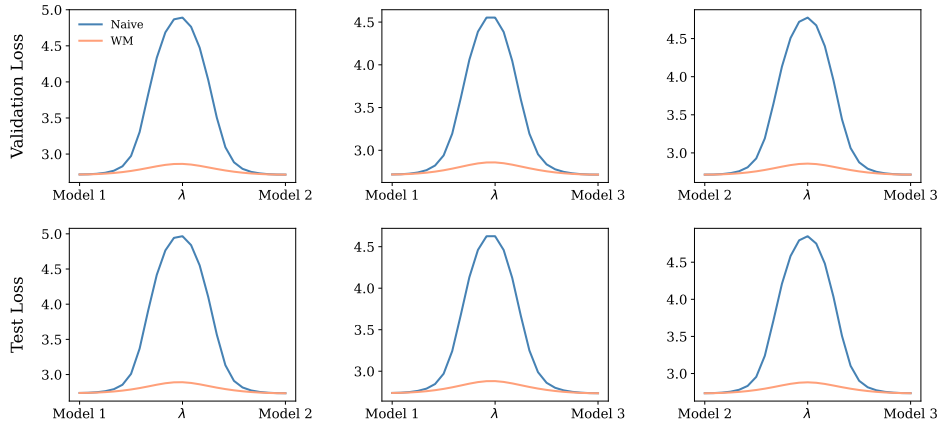


Figure 28: Linear Mode Connectivity for GPT2-MoE on Wikitext103 with 12 layers and 2 experts

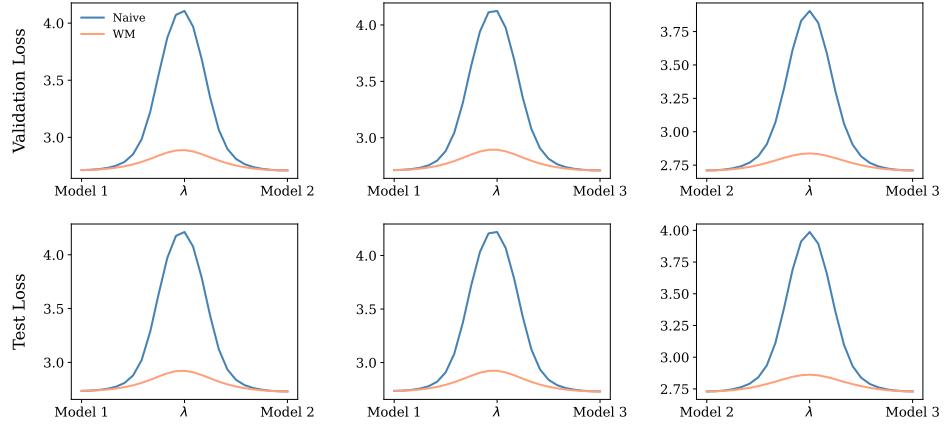


Figure 29: Linear Mode Connectivity for GPT2-MoE on Wikitext103 with 12 layers and 4 experts

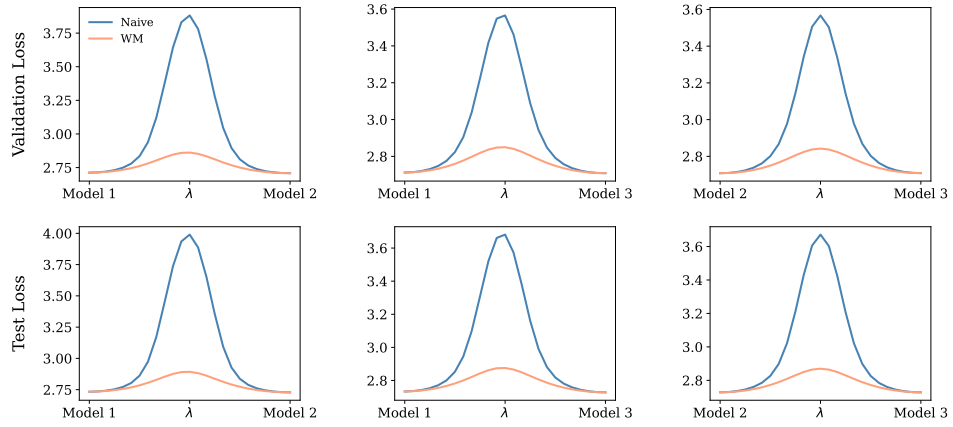


Figure 30: Linear Mode Connectivity for GPT2-MoE on Wikitext103 with 12 layers and 6 experts

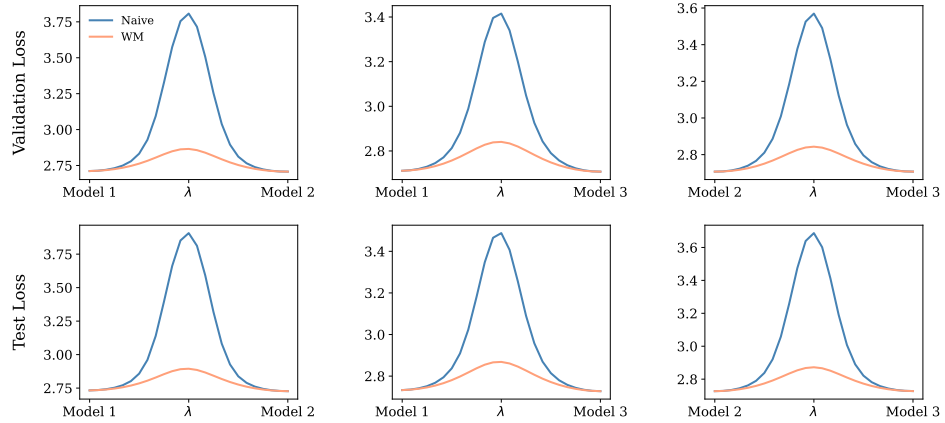


Figure 31: Linear Mode Connectivity for GPT2-MoE on Wikitext103 with 12 layers and 8 experts

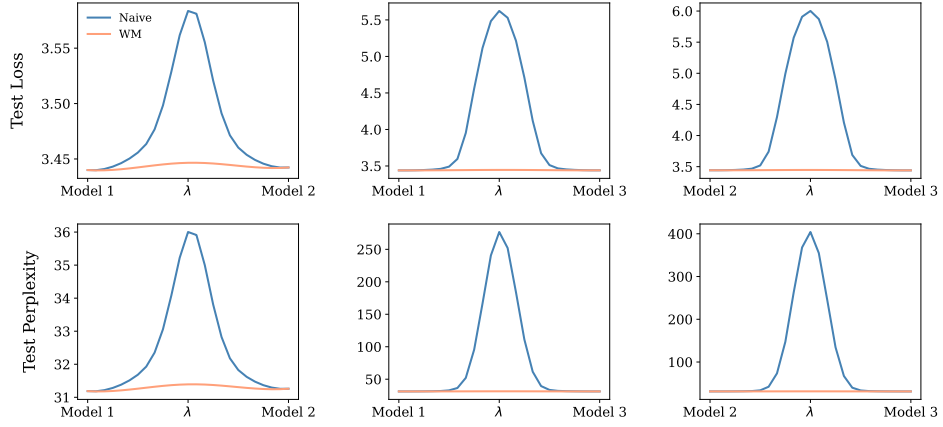


Figure 32: Linear Mode Connectivity for GPT2-MoE on One Billion Word with 12 layers and 2 experts

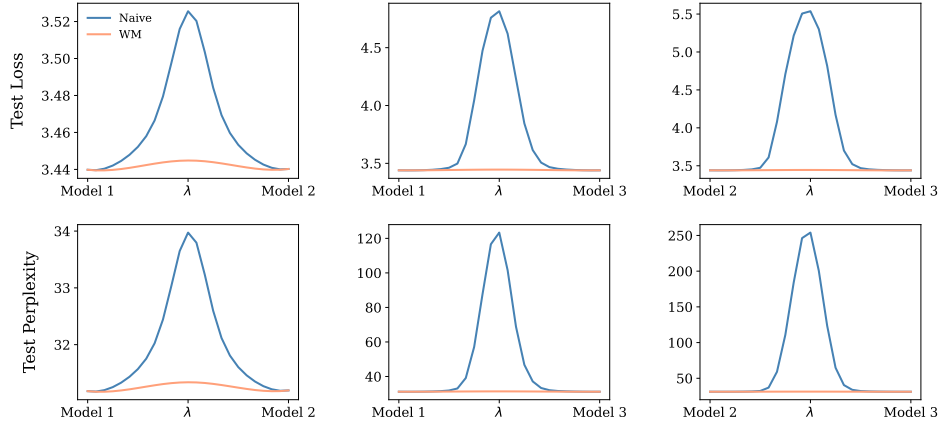


Figure 33: Linear Mode Connectivity for GPT2-MoE on One Billion Word with 12 layers and 4 experts

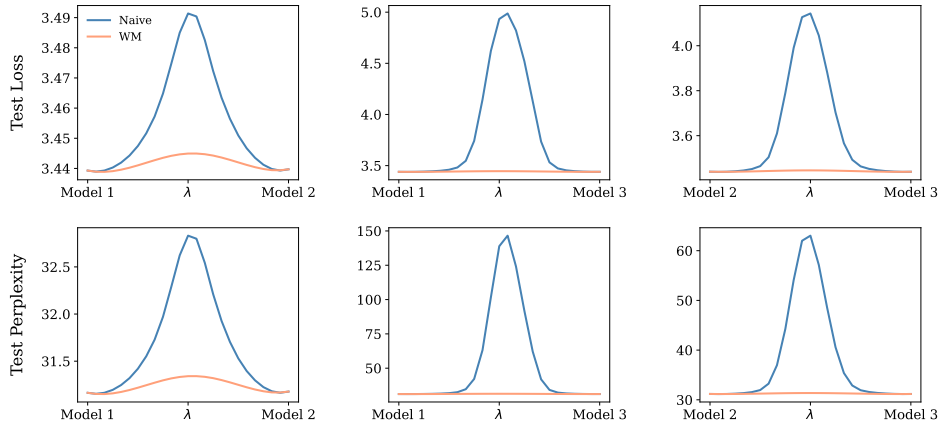


Figure 34: Linear Mode Connectivity for GPT2-MoE on One Billion Word with 12 layers and 6 experts

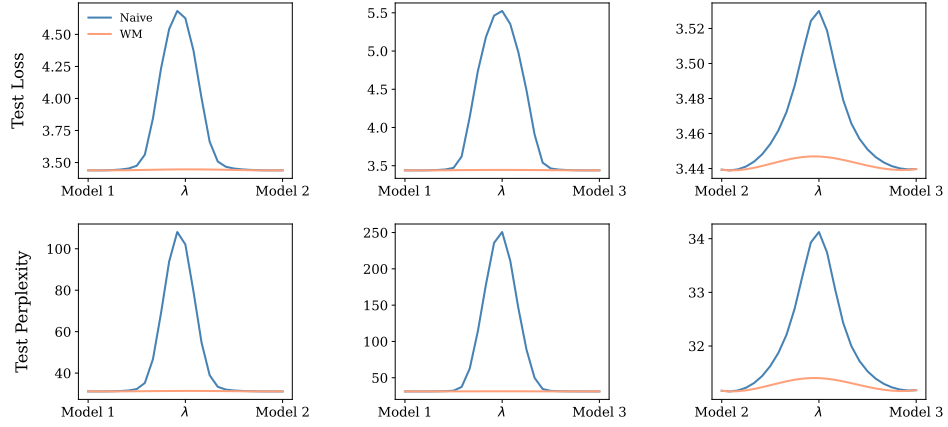


Figure 35: Linear Mode Connectivity for GPT2-MoE on One Billion Word with 12 layers and 8 experts

1299 G.1.2 Sparse Mixture-of-Experts

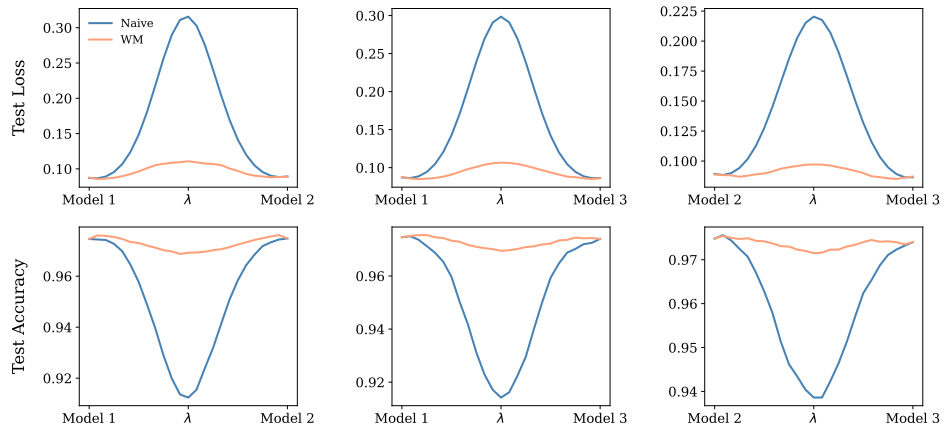


Figure 36: Linear Mode Connectivity for ViT-SMoE ($k = 2$) on MNIST with 1 layer and 4 experts

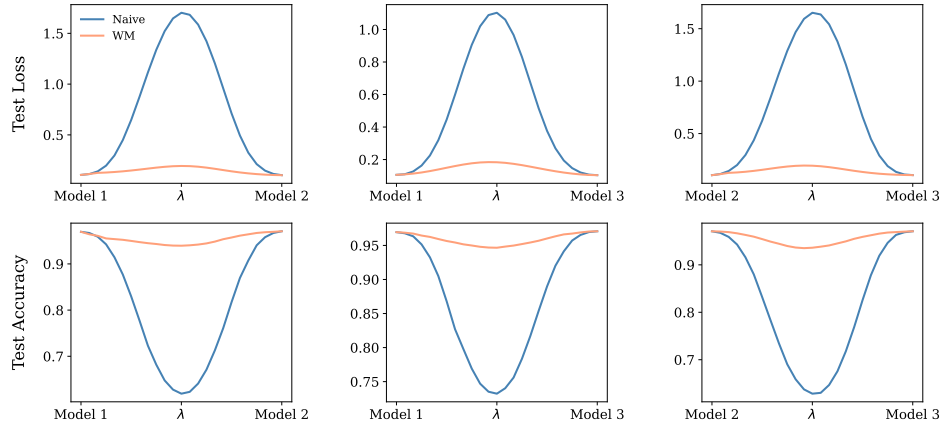


Figure 37: Linear Mode Connectivity for ViT-SMoE ($k = 2$) on MNIST with 2 layers and 4 experts

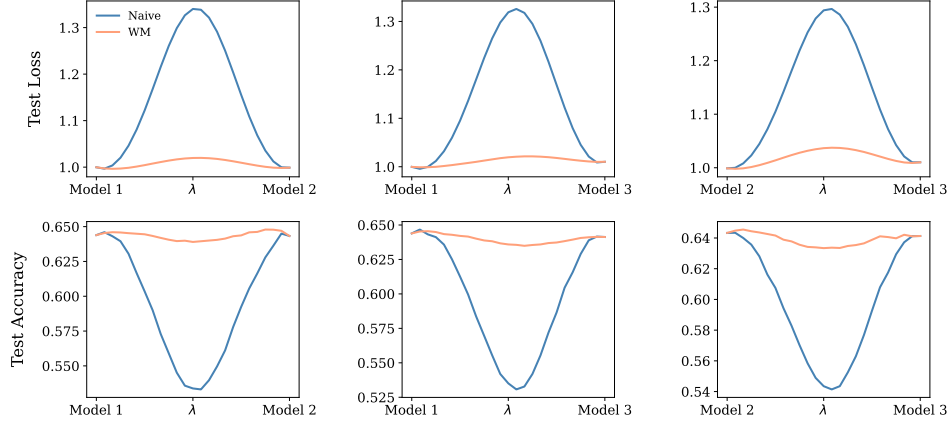


Figure 38: Linear Mode Connectivity for ViT-SMoE ($k = 2$) on CIFAR-10 with 2 layers and 4 experts

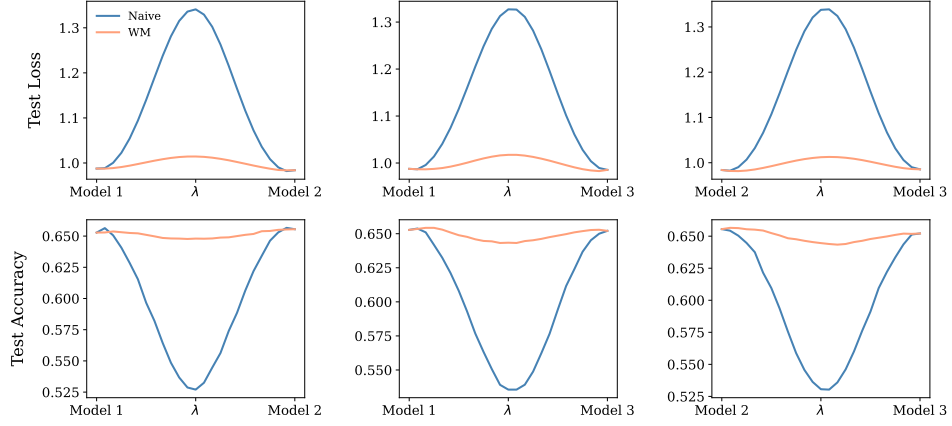


Figure 39: Linear Mode Connectivity for ViT-SMoE ($k = 2$) on CIFAR-10 with 2 layers and 8 experts

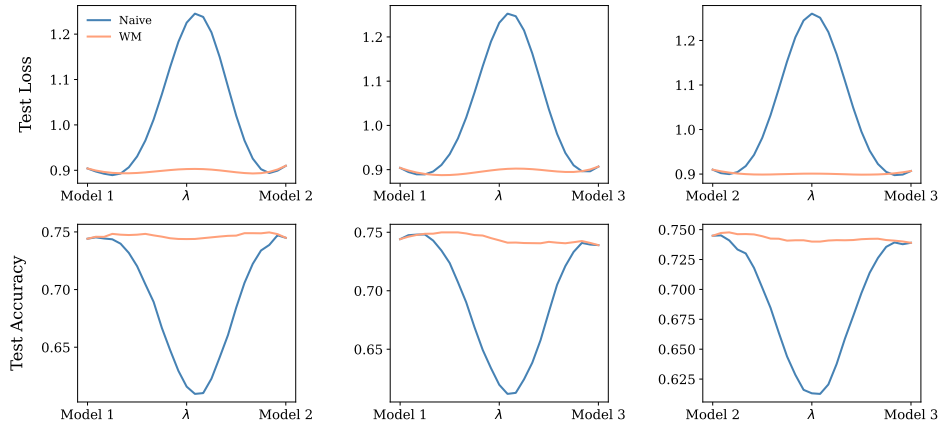


Figure 40: Linear Mode Connectivity for ViT-SMoE ($k = 2$) on CIFAR-10 with 6 layers and 4 experts

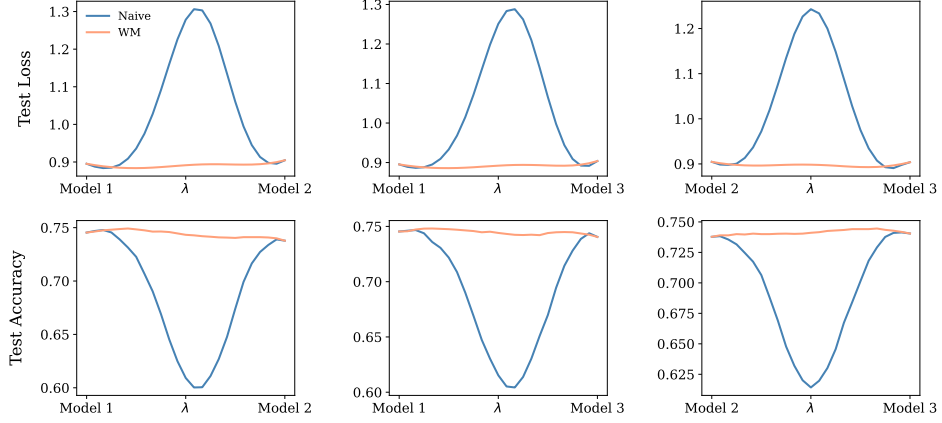


Figure 41: Linear Mode Connectivity for ViT-SMoE ($k = 2$) on CIFAR-10 with 6 layers and 8 experts

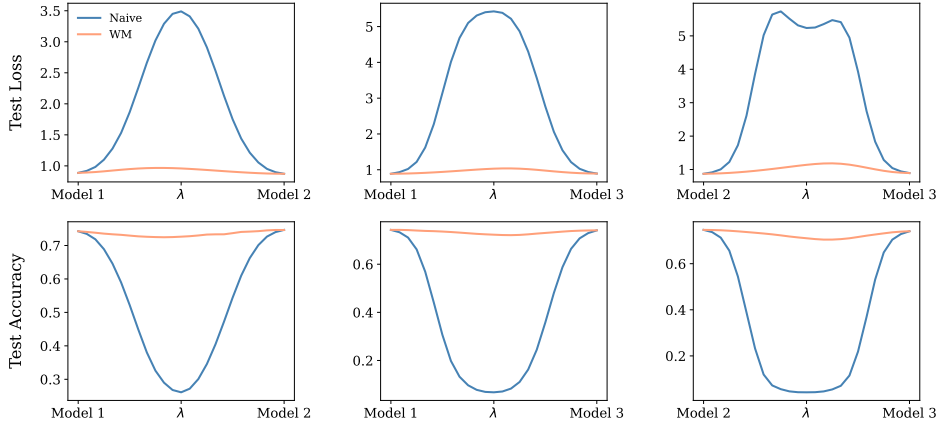


Figure 42: Linear Mode Connectivity for ViT-SMoE ($k = 2$) on CIFAR-100 with 6 layers and 4 experts

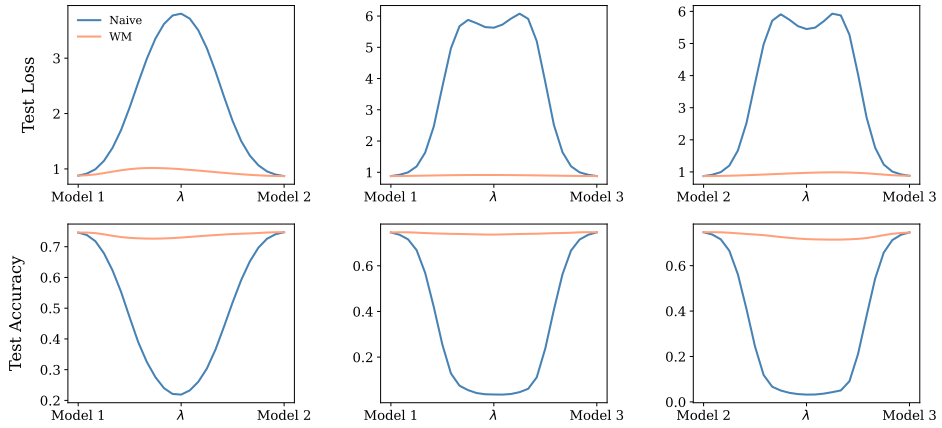


Figure 43: Linear Mode Connectivity for ViT-SMoE ($k = 2$) on CIFAR-100 with 6 layers and 8 experts

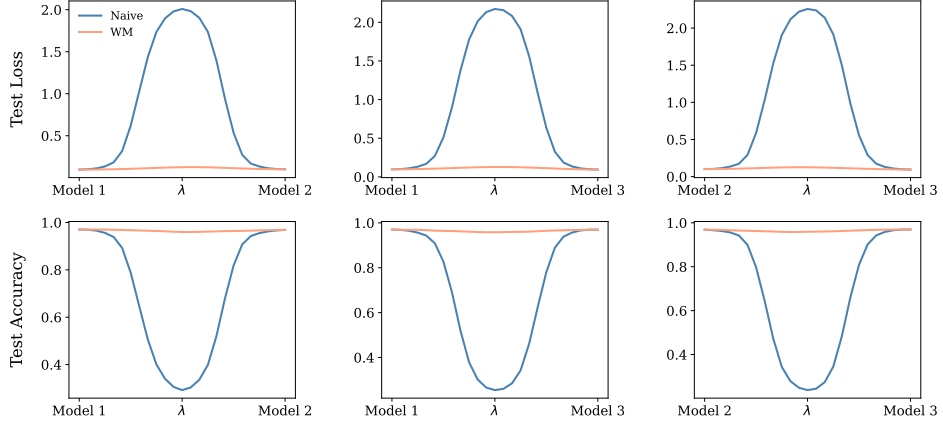


Figure 44: Linear Mode Connectivity for ViT-SMoE ($k = 2$) on ImageNet-21k→CIFAR-10 with 12 layers and 4 experts

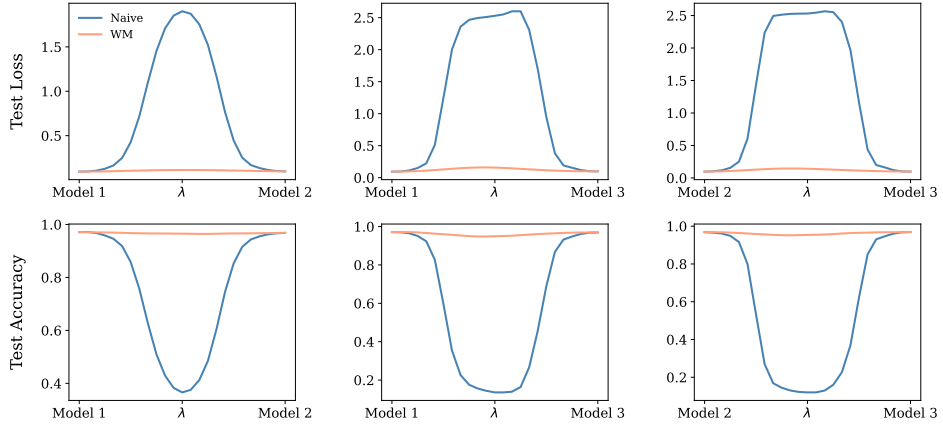


Figure 45: Linear Mode Connectivity for ViT-SMoE ($k = 2$) on ImageNet-21k→CIFAR-10 with 12 layers and 8 experts

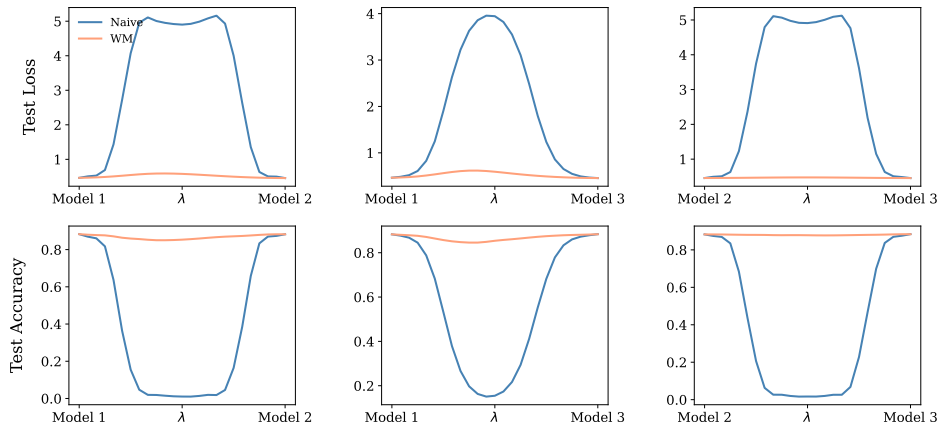


Figure 46: Linear Mode Connectivity for ViT-SMoE ($k = 2$) on ImageNet-21k→CIFAR-100 with 12 layers and 4 experts

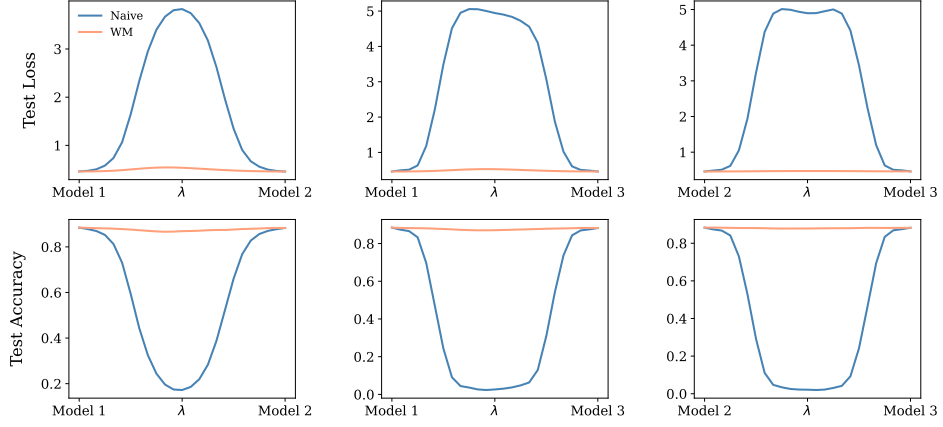


Figure 47: Linear Mode Connectivity for ViT-SMoE ($k = 2$) on ImageNet-21k→CIFAR-100 with 12 layers and 8 experts

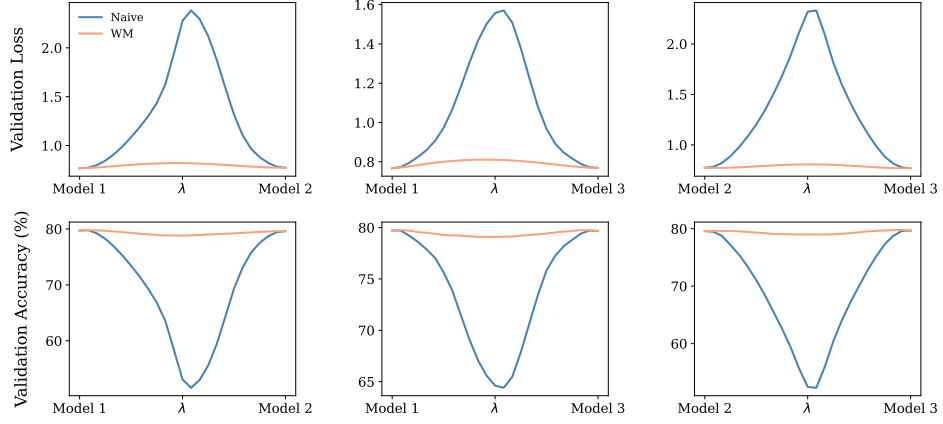


Figure 48: Linear Mode Connectivity for ViT-SMoE ($k = 2$) on ImageNet-1k with 12 layers and 4 experts

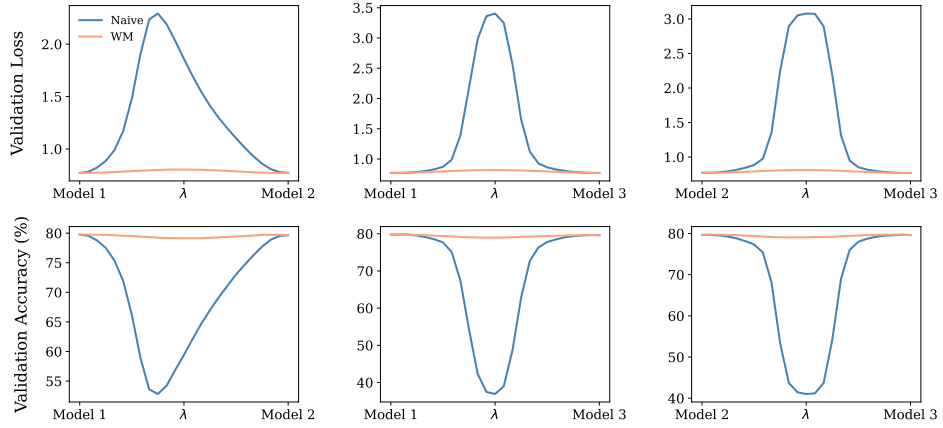


Figure 49: Linear Mode Connectivity for ViT-SMoE ($k = 2$) on ImageNet-1k with 12 layers and 8 experts

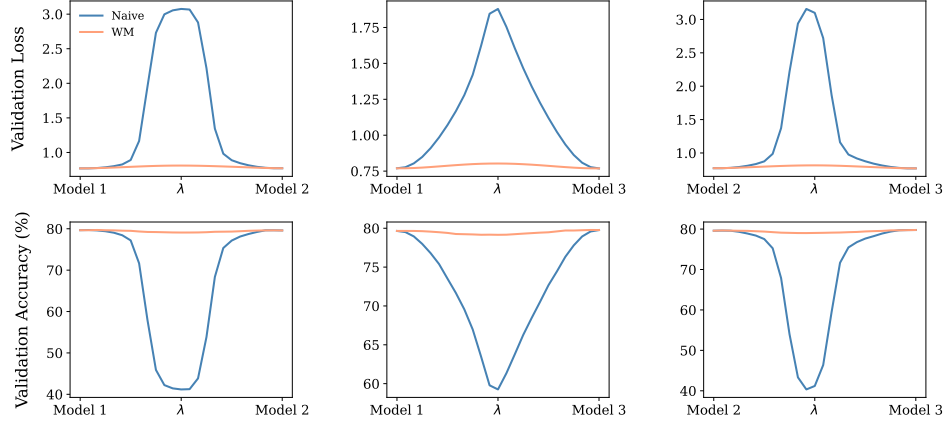


Figure 50: Linear Mode Connectivity for ViT-SMoE ($k = 2$) on ImageNet-1k with 12 layers and 16 experts

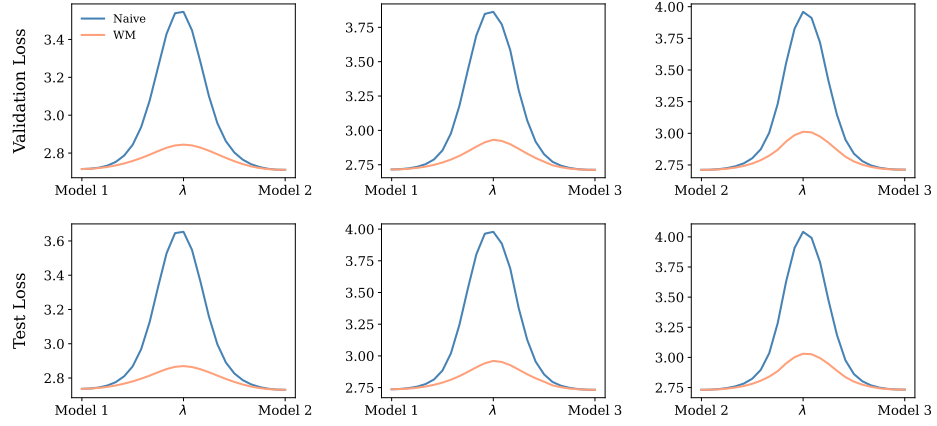


Figure 51: Linear Mode Connectivity for GPT2-SMoE ($k = 2$) on Wikitext103 with 12 layers and 4 experts

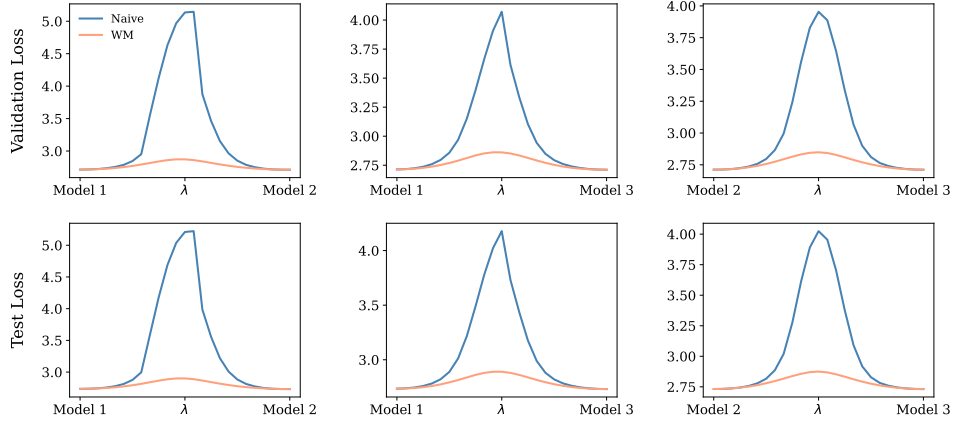


Figure 52: Linear Mode Connectivity for GPT2-SMoE ($k = 2$) on Wikitext103 with 12 layers and 8 experts

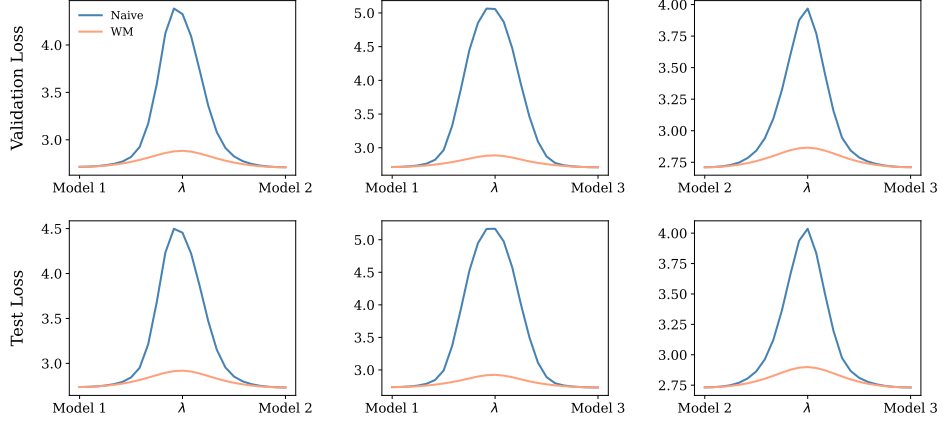


Figure 53: Linear Mode Connectivity for GPT2-SMoE ($k = 2$) on Wikitext103 with 12 layers 16 experts

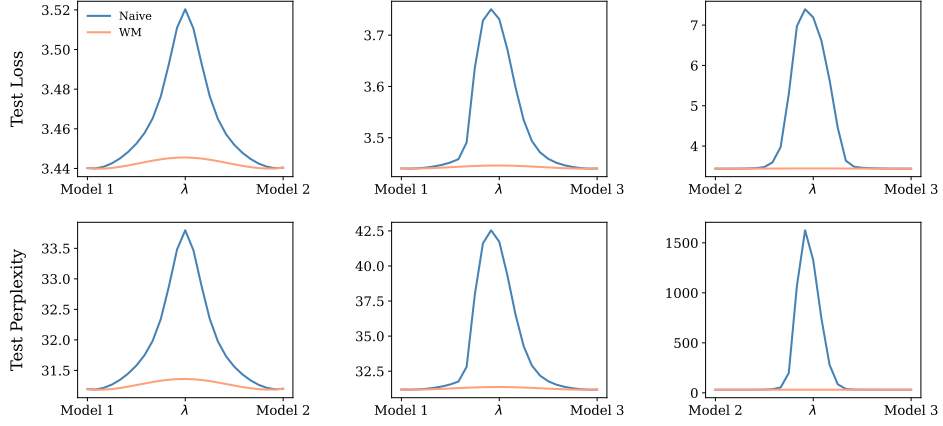


Figure 54: Linear Mode Connectivity for GPT2-SMoE ($k = 2$) on One Billion Word with 12 layers and 4 experts

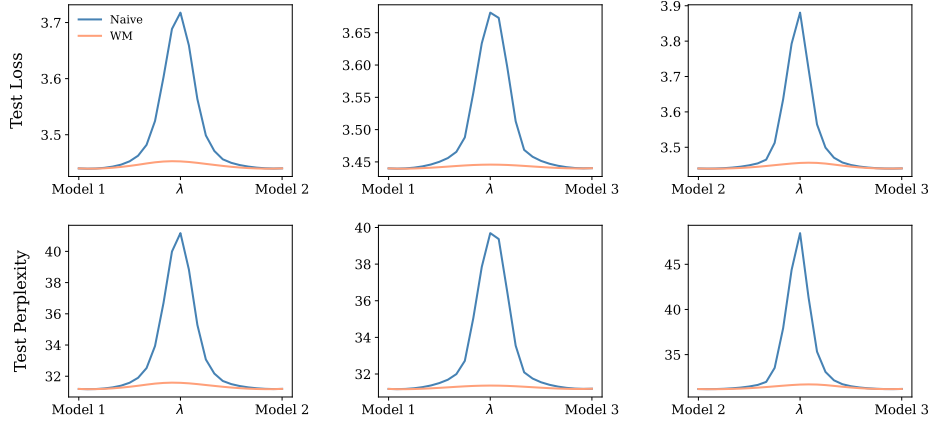


Figure 55: Linear Mode Connectivity for GPT2-SMoE ($k = 2$) on One Billion Word with 12 layers and 8 experts

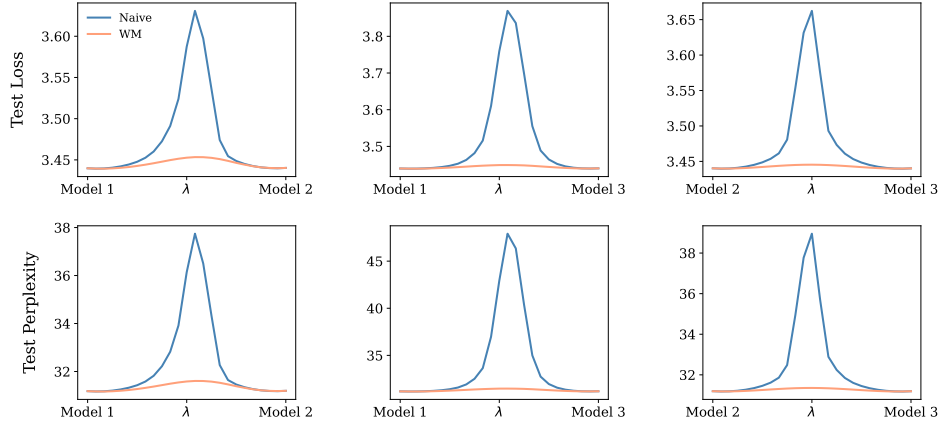


Figure 56: Linear Mode Connectivity for GPT2-SMoE ($k = 2$) on One Billion Word with 12 layers and 16 experts

1300 G.1.3 DeepSeek Mixture-of-Experts

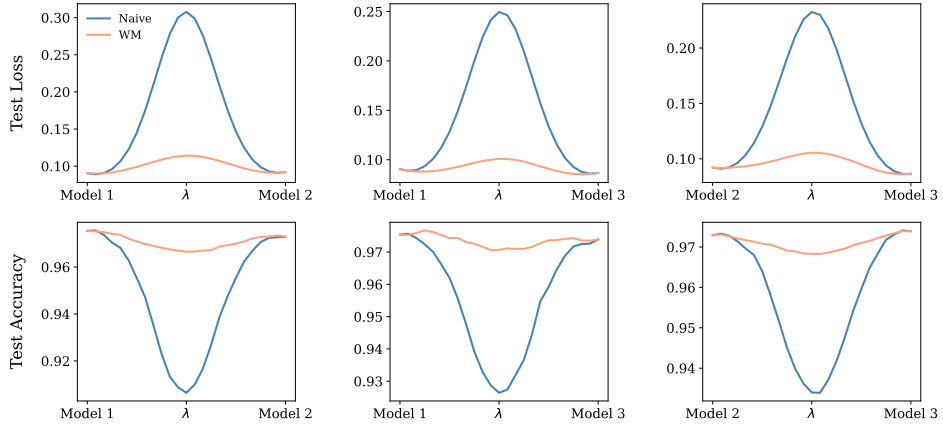


Figure 57: Linear Mode Connectivity for ViT-DeepSeekMoE ($k = 2, s = 1$) on MNIST with 1 layer and 4 experts

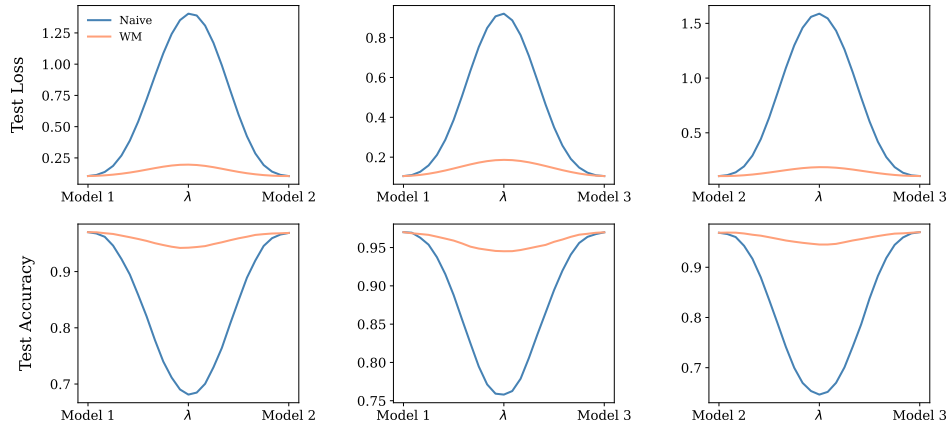


Figure 58: Linear Mode Connectivity for ViT-DeepSeekMoE ($k = 2, s = 1$) on MNIST with 2 layers and 4 experts

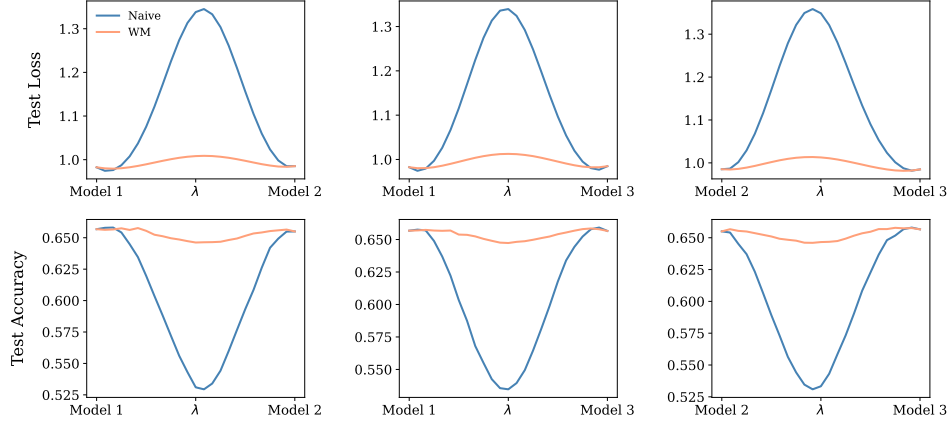


Figure 59: Linear Mode Connectivity for ViT-DeepSeekMoE ($k = 2, s = 1$) on CIFAR-10 with 2 layers and 4 experts

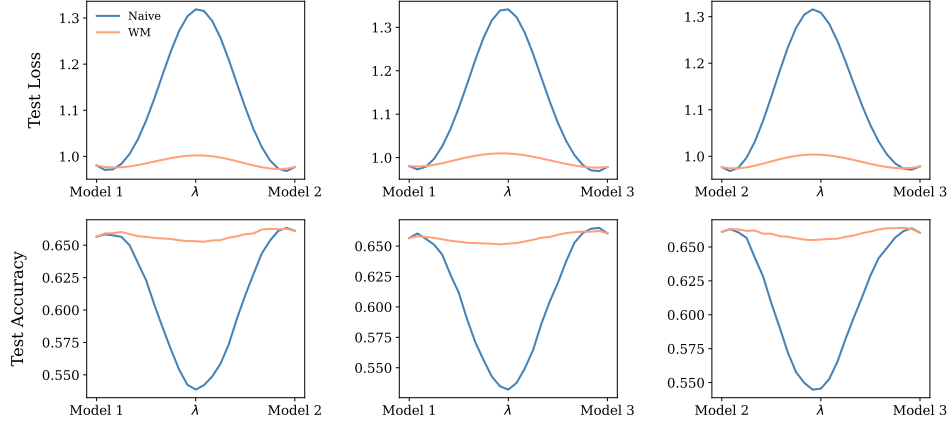


Figure 60: Linear Mode Connectivity for ViT-DeepSeekMoE ($k = 2, s = 1$) on CIFAR-10 with 2 layers and 8 experts

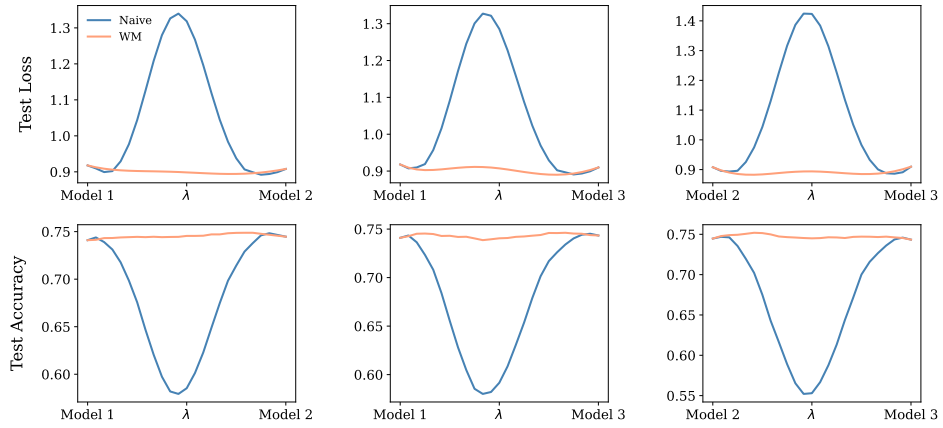


Figure 61: Linear Mode Connectivity for ViT-DeepSeekMoE ($k = 2, s = 1$) on CIFAR-10 with 6 layers and 4 experts

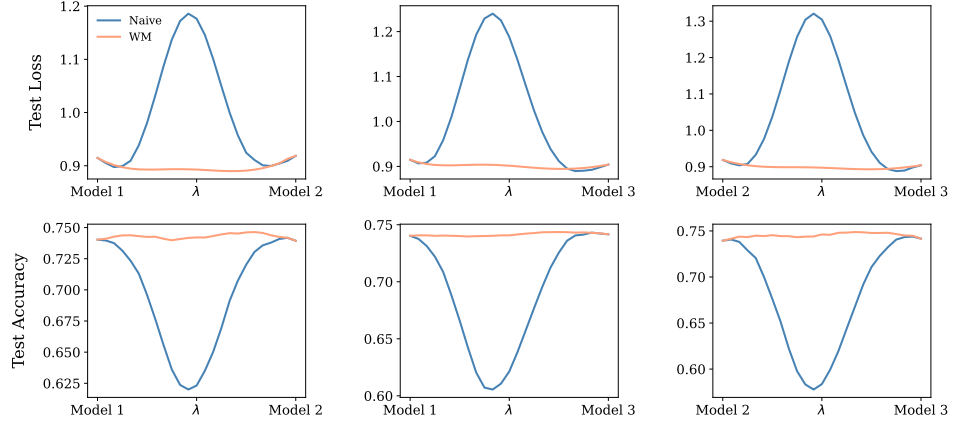


Figure 62: Linear Mode Connectivity for ViT-DeepSeekMoE ($k = 2, s = 1$) on CIFAR-10 with 6 layers and 8 experts

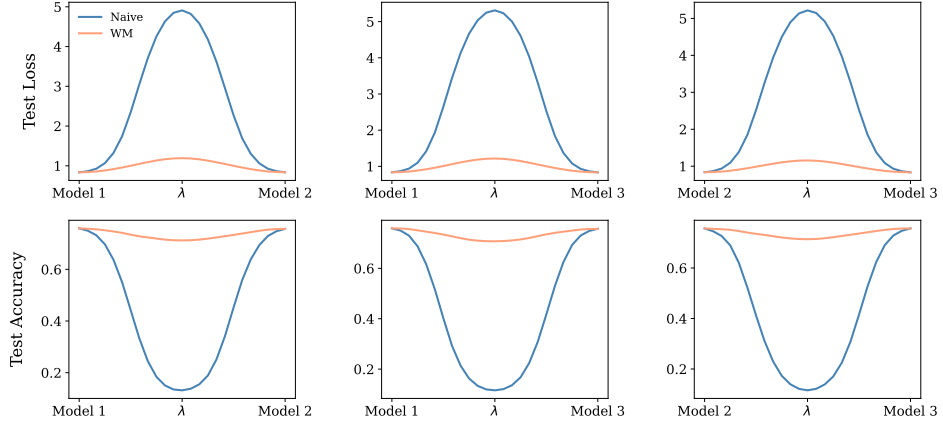


Figure 63: Linear Mode Connectivity for ViT-DeepSeekMoE ($k = 2, s = 1$) on CIFAR-100 with 6 layers and 4 experts

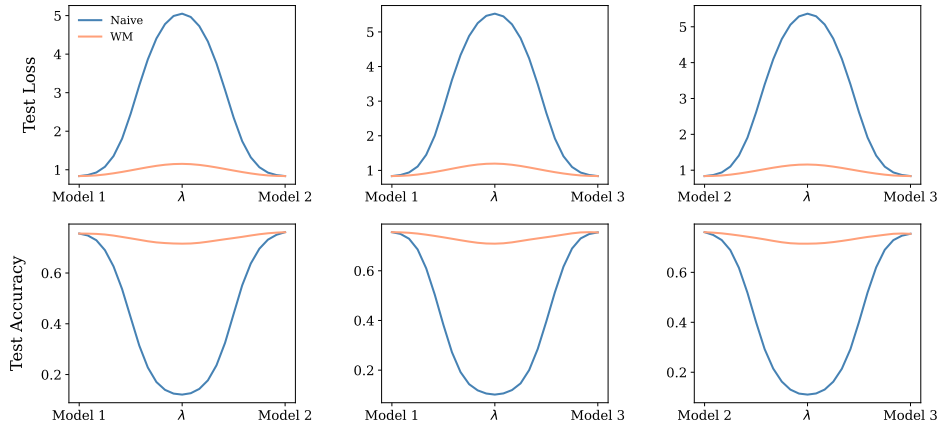


Figure 64: Linear Mode Connectivity for ViT-DeepSeekMoE ($k = 2, s = 1$) on CIFAR-100 with 6 layers and 8 experts

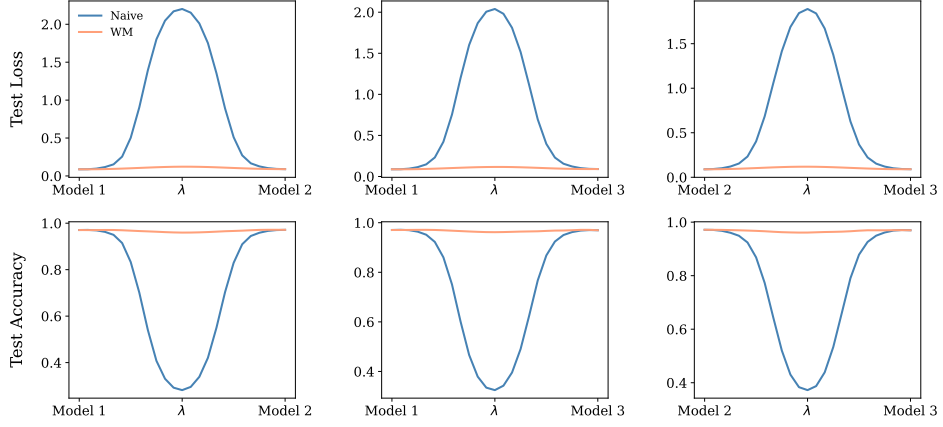


Figure 66: Linear Mode Connectivity for ViT-DeepSeekMoE ($k = 2, s = 1$) on ImageNet-21k→CIFAR-10 with 12 layers and 8 experts

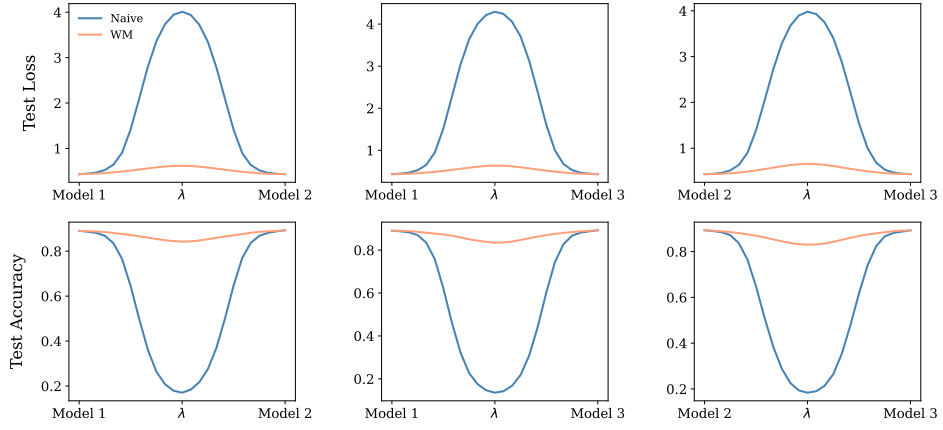


Figure 67: Linear Mode Connectivity for ViT-DeepSeekMoE ($k = 2, s = 1$) on ImageNet-21k→CIFAR-100 with 12 layers and 4 experts

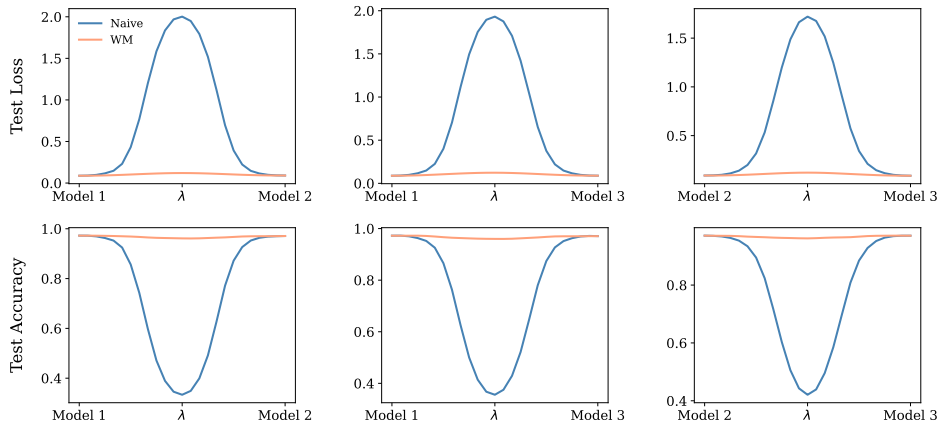


Figure 65: Linear Mode Connectivity for ViT-DeepSeekMoE ($k = 2, s = 1$) on ImageNet-21k→CIFAR-10 with 12 layers and 4 experts

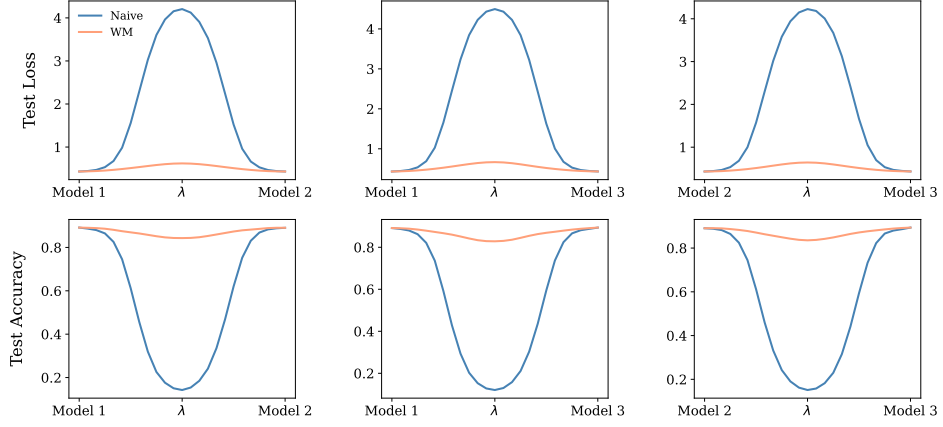


Figure 68: Linear Mode Connectivity for ViT-DeepSeekMoE ($k = 2, s = 1$) on ImageNet-21k→CIFAR-100 with 12 layers and 8 experts

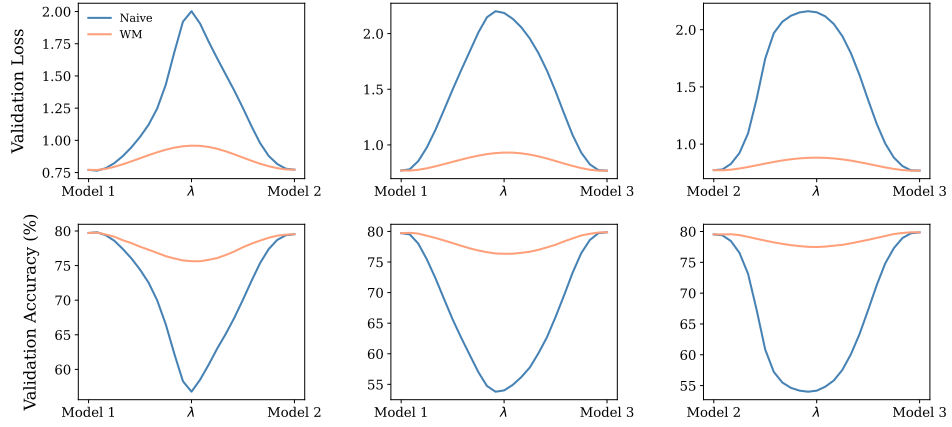


Figure 69: Linear Mode Connectivity for ViT-DeepSeekMoE ($k = 2, s = 1$) on ImageNet-1k with 12 layers and 4 experts

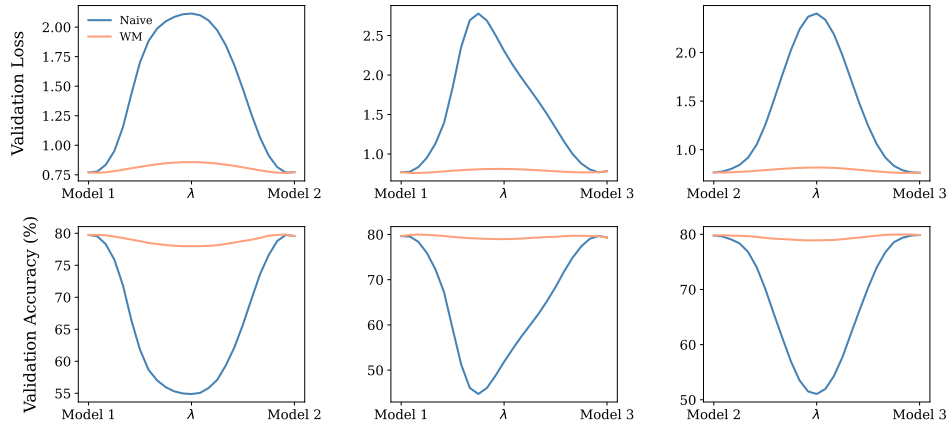


Figure 70: Linear Mode Connectivity for ViT-DeepSeekMoE ($k = 2, s = 1$) on ImageNet-1k with 12 layers and 8 experts

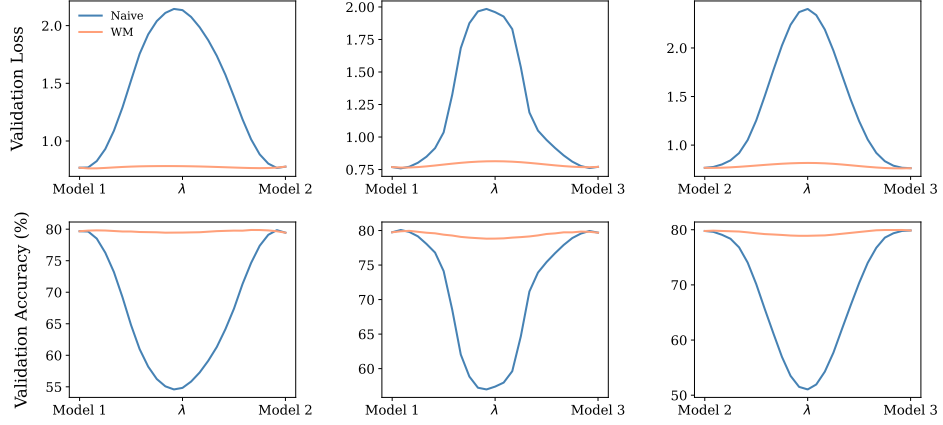


Figure 71: Linear Mode Connectivity for ViT-DeepSeekMoE ($k = 2, s = 1$) on ImageNet-1k with 12 layers and 16 experts

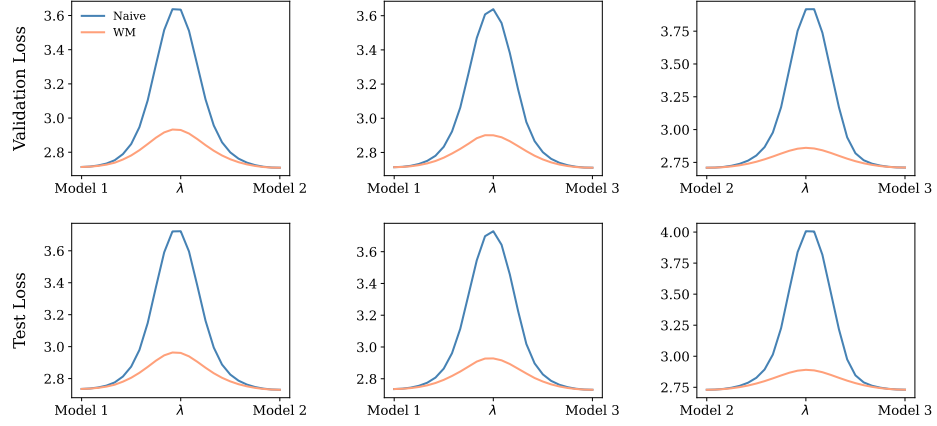


Figure 72: Linear Mode Connectivity for ViT-DeepSeekMoE ($k = 2, s = 1$) on Wikitext103 with 12 layers and 4 experts

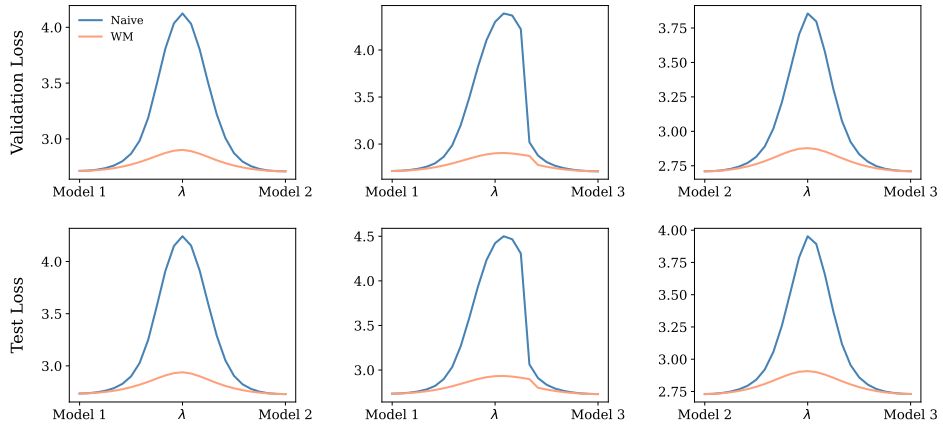


Figure 73: Linear Mode Connectivity for ViT-DeepSeekMoE ($k = 2, s = 1$) on Wikitext103 with 12 layers and 8 experts

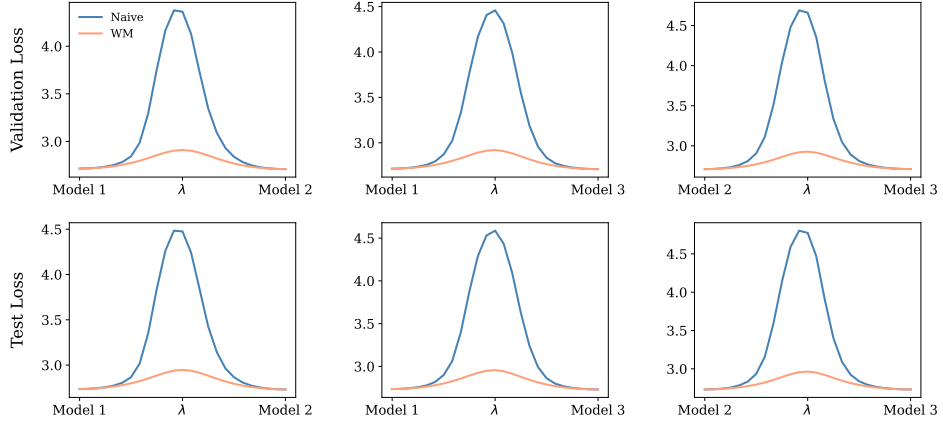


Figure 74: Linear Mode Connectivity for ViT-DeepSeekMoE ($k = 2, s = 1$) on Wikitext103 with 12 layers and 16 experts

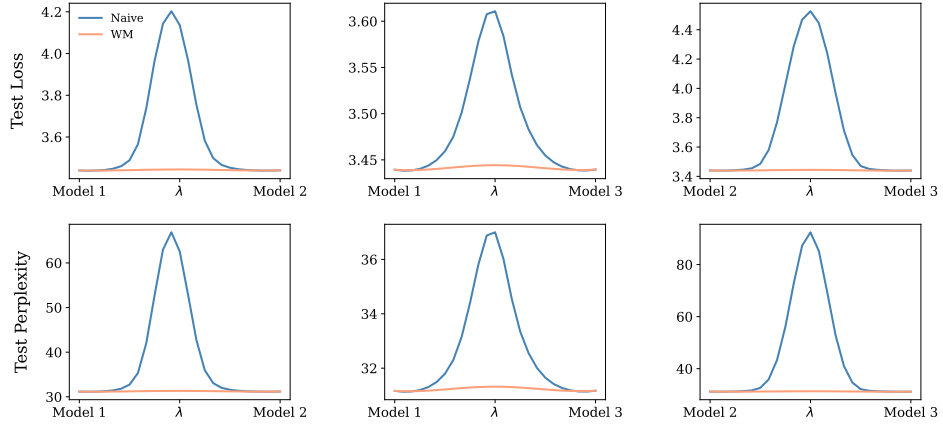


Figure 75: Linear Mode Connectivity for ViT-DeepSeekMoE ($k = 2, s = 1$) on One Billion Word with 12 layers and 4 experts

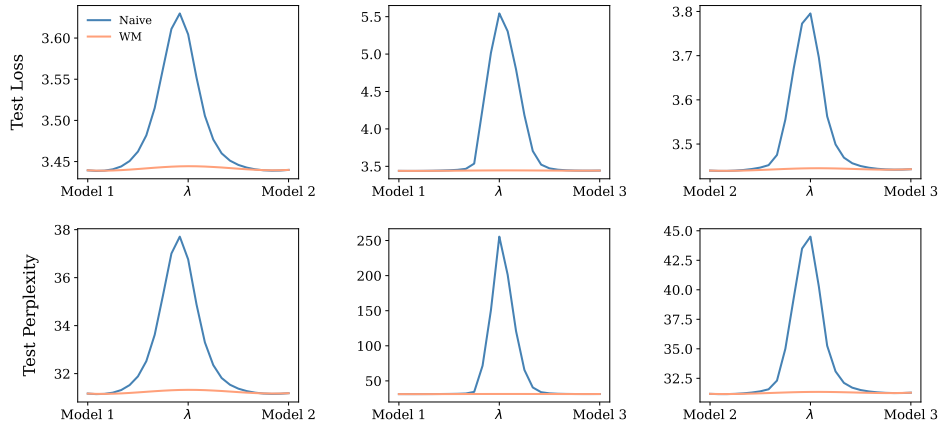


Figure 76: Linear Mode Connectivity for ViT-DeepSeekMoE ($k = 2, s = 1$) on One Billion Word with 12 layers and 8 experts

Table 6: In all configurations, the MLP component of the *last* Transformer layer is replaced by an MoE component.

Method	Dataset	Number of layers	Number of experts	Figure
MoE	CIFAR-10	6	[2, 4]	[78, 79]
	ImageNet-21k→CIFAR-10	12	[2, 4]	[80, 81]
	ImageNet-21k→CIFAR-100	12	[2, 4]	[82, 83]
SMoE ($k = 2$)	CIFAR-10	6	[4, 8]	[84, 85]
	ImageNet-21k→CIFAR-10	12	[4, 8]	[86, 87]
	ImageNet-21k→CIFAR-100	12	[4, 8]	[88, 89]
DeepSeekMoE ($k = 2, s = 1$)	CIFAR-10	6	[4, 8]	[90, 91]
	ImageNet-21k→CIFAR-10	12	[4, 8]	[92, 93]
	ImageNet-21k→CIFAR-100	12	[4, 8]	[94, 95]

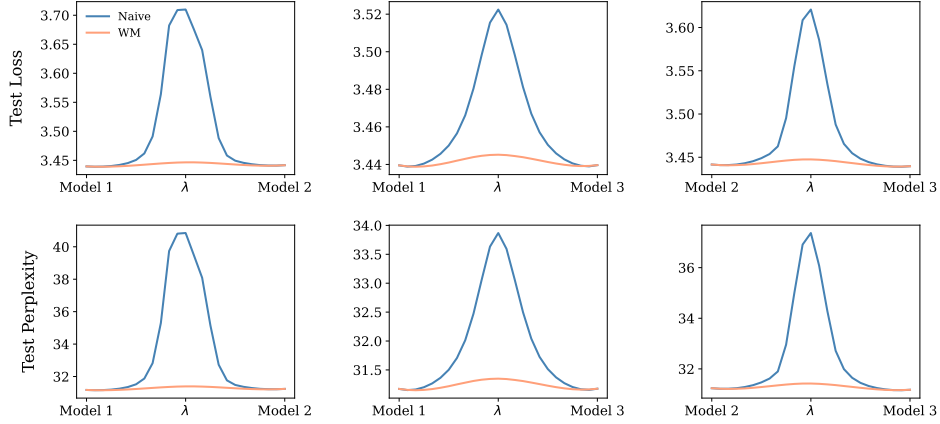


Figure 77: Linear Mode Connectivity for ViT-DeepSeekMoE ($k = 2, s = 1$) on One Billion Word with 12 layers and 16 experts

G.2 Linear Mode Connectivity Analysis: Last Layer

We investigate Linear Mode Connectivity (LMC) in settings where the feed-forward network (FFN) of the *last* Transformer layer is replaced with a Mixture-of-Experts (MoE) module. This extends our earlier analysis of *first*-layer substitutions by examining connectivity at the terminal depth of the architecture, offering a complementary view on model plasticity and the modular structure of expert-based designs.

Following the procedure outlined in Section 6.1, all pretrained weights were frozen, and only the MoE module in the *last* layer was fine-tuned. For each configuration, three independent fine-tuning runs were conducted with different random seeds. LMC was assessed by evaluating the loss along linear interpolation paths between pairs of fine-tuned models, providing insight into the connectivity and overlap of their respective solution basins.

Table 6 summarizes the experimental settings for *last*-layer FFN replacement across dense MoE, SMoE, and DeepSeekMoE variants. Unless otherwise stated, all figure references correspond to *last*-layer experiments.

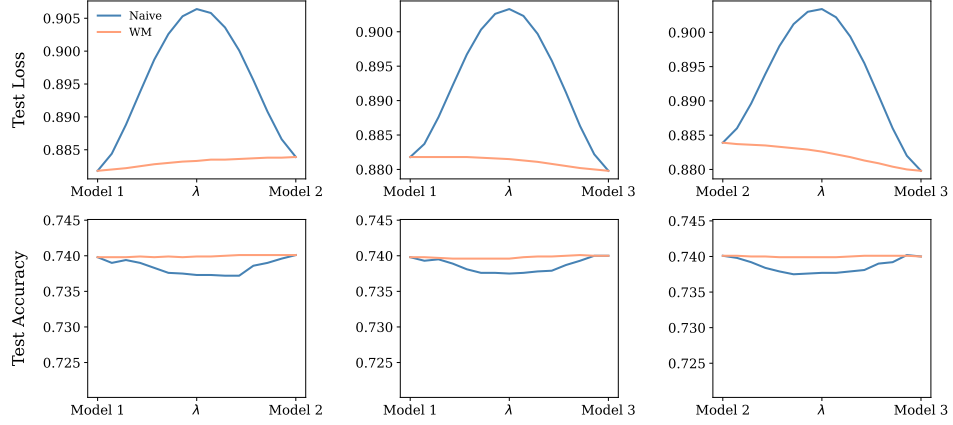


Figure 78: Linear Mode Connectivity for ViT-MoE on CIFAR-10 with 6 layers and 2 experts.

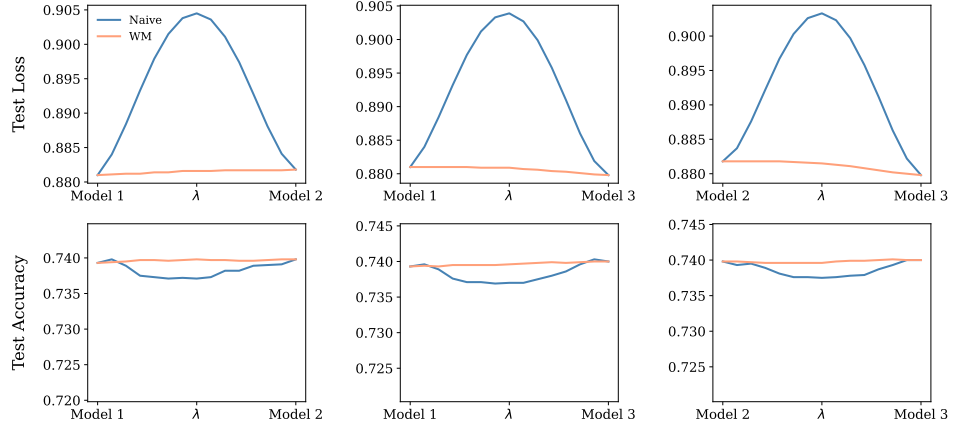


Figure 79: Linear Mode Connectivity for ViT-MoE on CIFAR-10 with 6 layers and 4 experts.

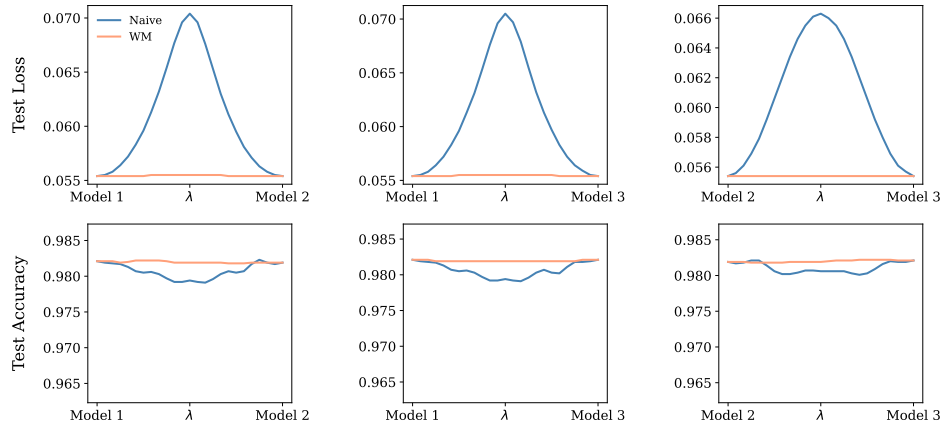


Figure 80: Linear Mode Connectivity for ViT-MoE on ImageNet-21k→CIFAR-10 with 12 layers and 2 experts.

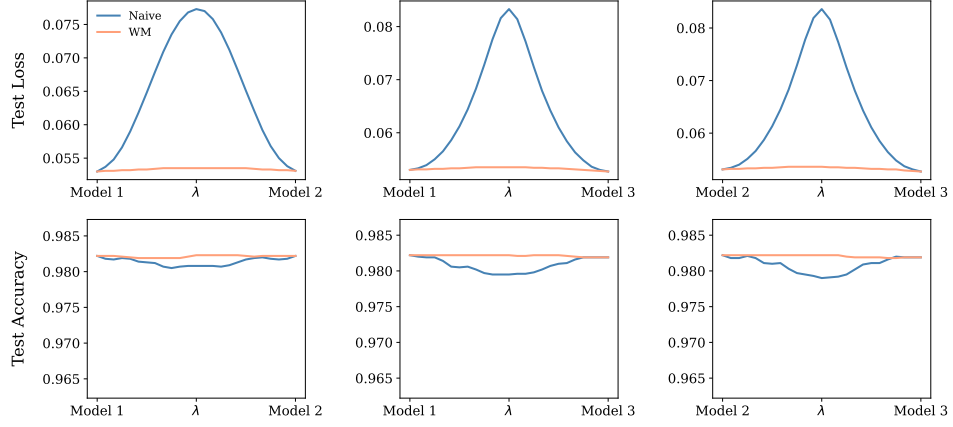


Figure 81: Linear Mode Connectivity for ViT-MoE on on ImageNet-21k→CIFAR-10 with 12 layers and 4 experts.

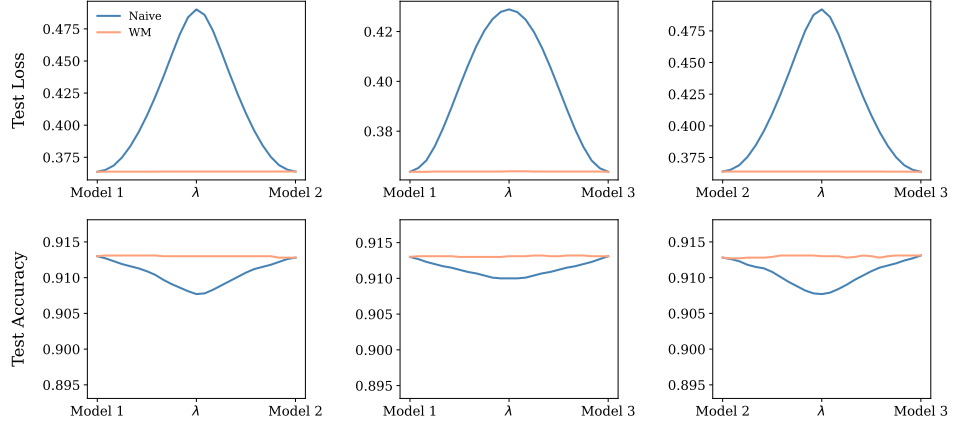


Figure 82: Linear Mode Connectivity for ViT-MoE on ImageNet-21k→CIFAR-100 with 12 layers and 2 experts.

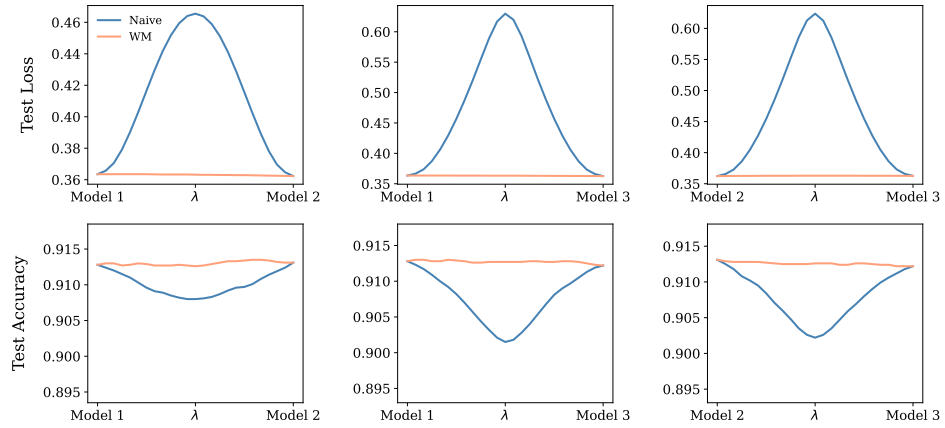


Figure 83: Linear Mode Connectivity for ViT-MoE on ImageNet-21k→CIFAR-100 with 12 layers and 4 experts.

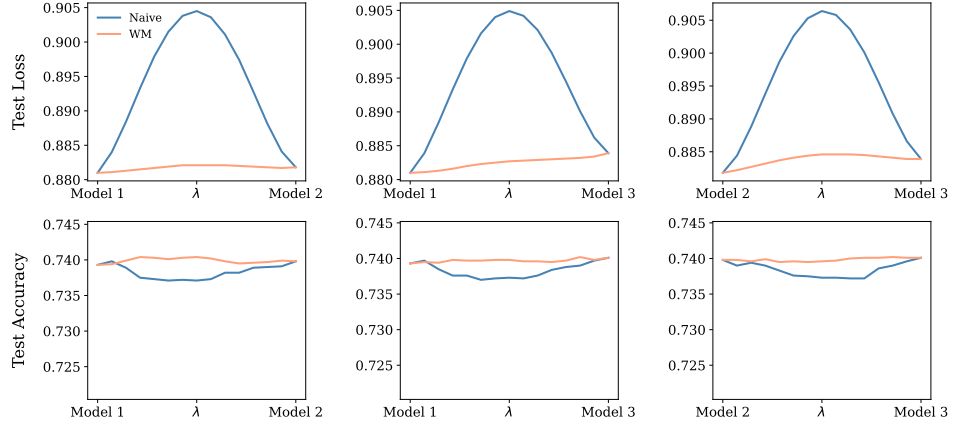


Figure 84: Linear Mode Connectivity for ViT-SMoE ($k = 2$) on CIFAR-10 with 6 layers and 4 experts.

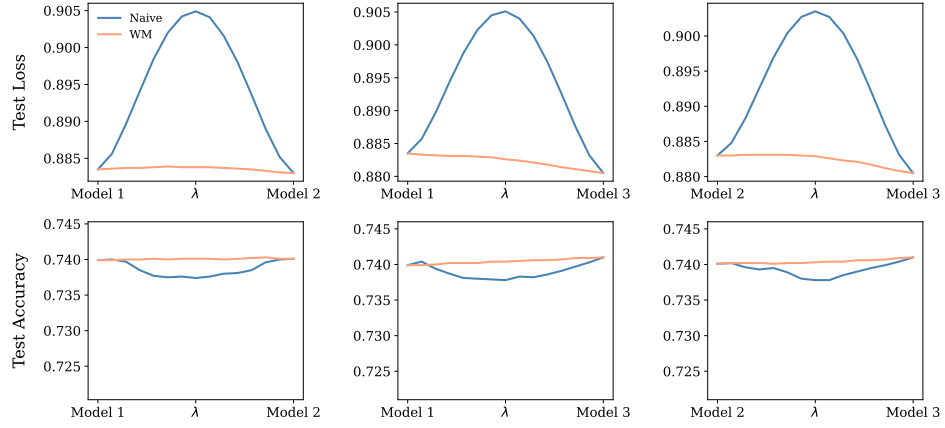


Figure 85: Linear Mode Connectivity for ViT-SMoE ($k = 2$) on CIFAR-10 with 6 layers and 8 experts.

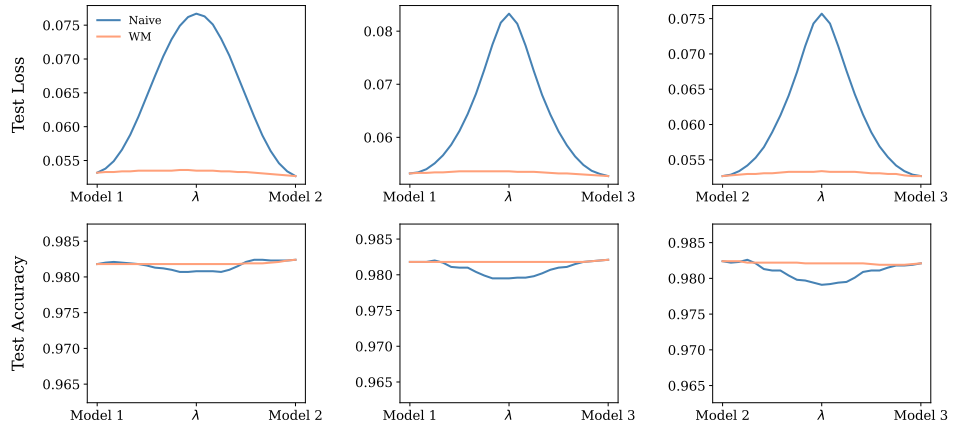


Figure 86: Linear Mode Connectivity for ViT-SMoE ($k = 2$) on ImageNet-21k→CIFAR-10 with 12 layers and 4 experts.

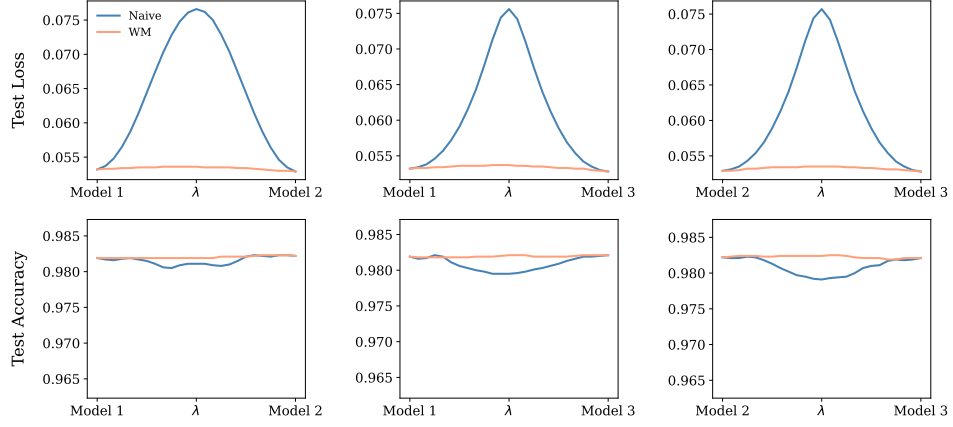


Figure 87: Linear Mode Connectivity for ViT-SMoE ($k = 2$) on ImageNet-21k→CIFAR-10 with 12 layers and 8 experts.

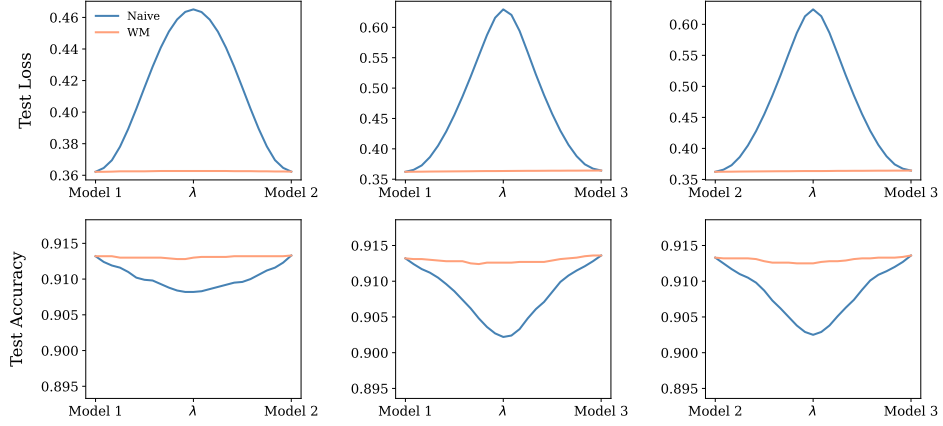


Figure 88: Linear Mode Connectivity for ViT-SMoE ($k = 2$) on ImageNet-21k→CIFAR-100 with 12 layers and 4 experts.

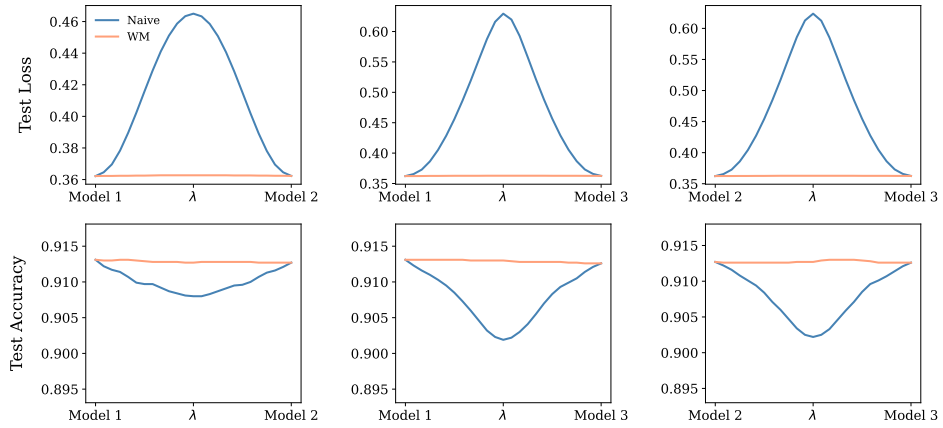


Figure 89: Linear Mode Connectivity for ViT-SMoE ($k = 2$) on ImageNet-21k→CIFAR-100 with 12 layers and 8 experts.

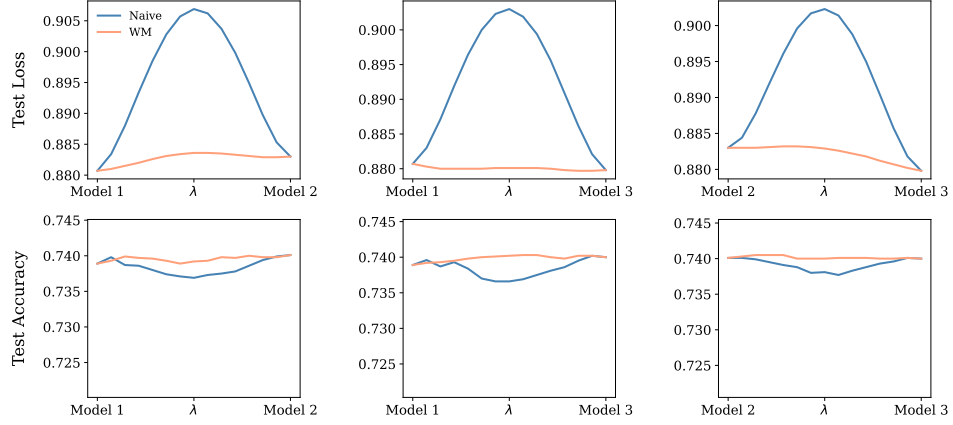


Figure 90: Linear Mode Connectivity for ViT-DeepSeekMoE ($k = 2, s = 1$) on CIFAR-10 with 6 layers and 4 experts.

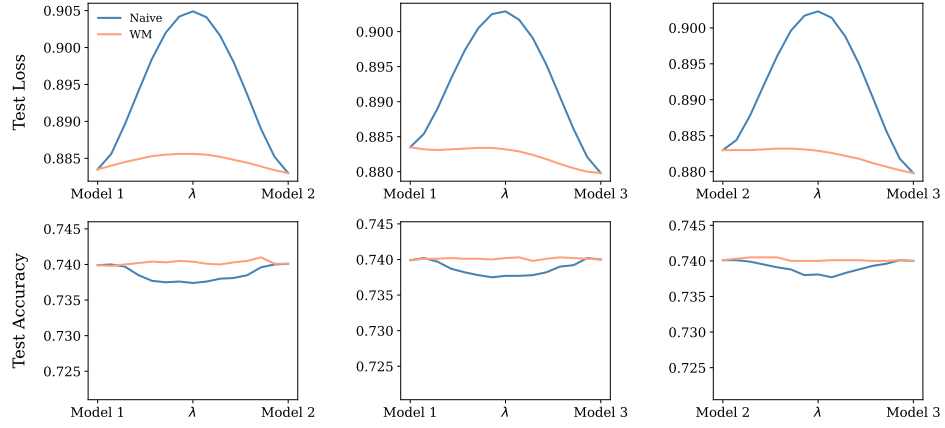


Figure 91: Linear Mode Connectivity for ViT-DeepSeekMoE ($k = 2, s = 1$) on CIFAR-10 with 6 layers and 8 experts.

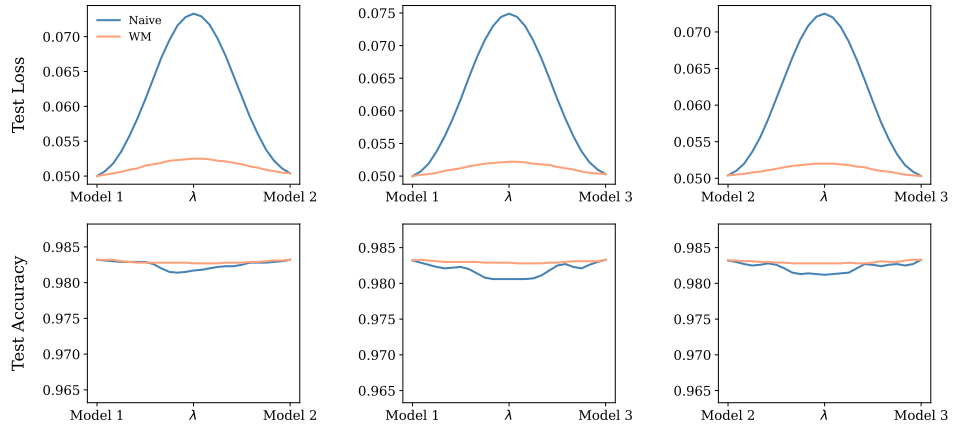


Figure 92: Linear Mode Connectivity for ViT-DeepSeekMoE ($k = 2, s = 1$) on ImageNet-21k \rightarrow CIFAR-10 with 12 layers and 4 experts.

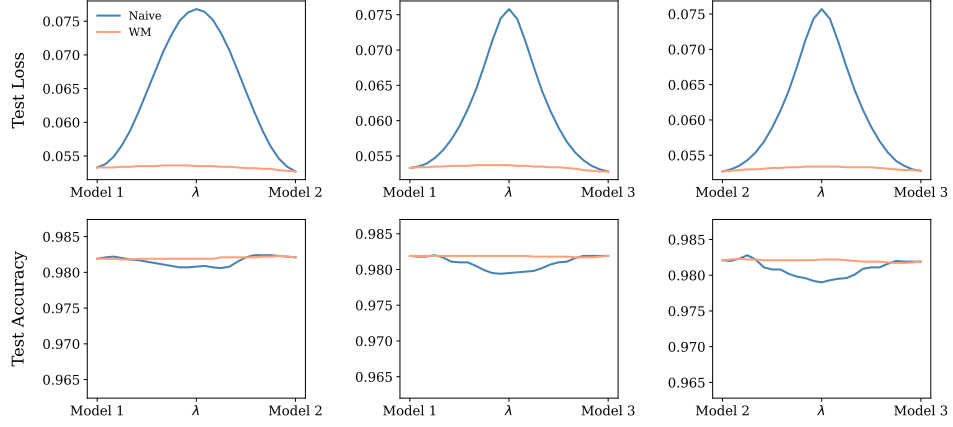


Figure 93: Linear Mode Connectivity for ViT-DeepSeekMoE ($k = 2, s = 1$) on ImageNet-21k→CIFAR-10 with 12 layers and 8 experts.

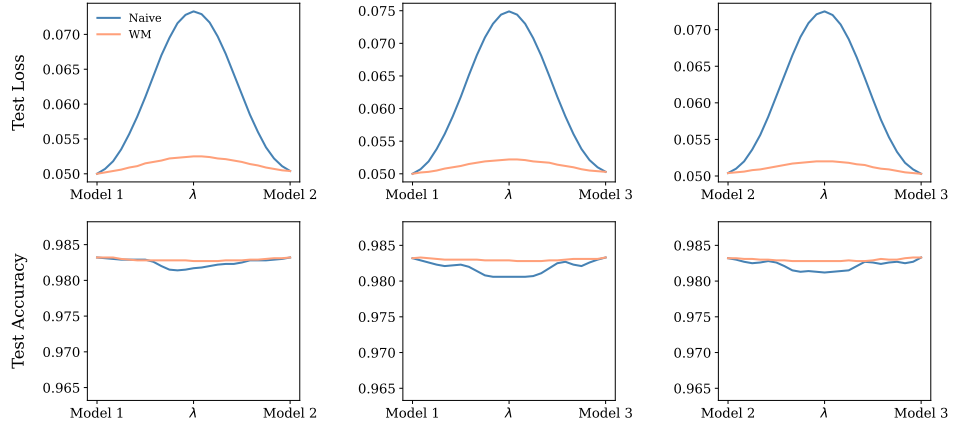


Figure 94: Linear Mode Connectivity for ViT-DeepSeekMoE ($k = 2, s = 1$) on ImageNet-21k→CIFAR-100 with 12 layers and 4 experts.

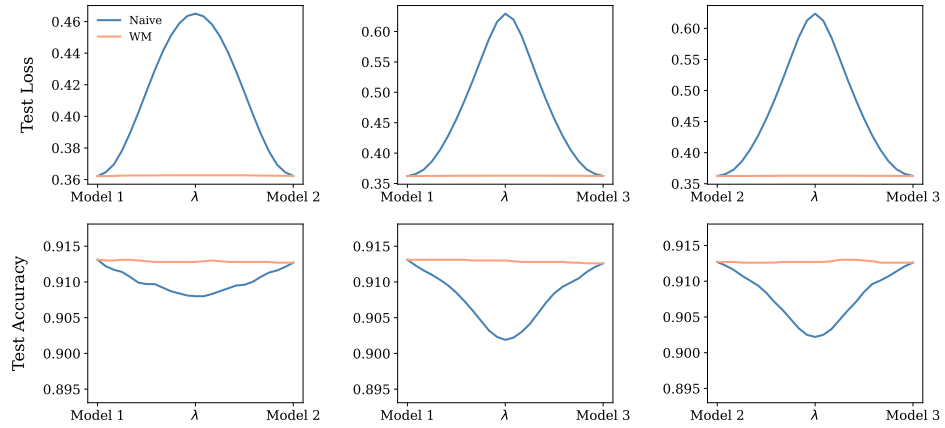


Figure 95: Linear Mode Connectivity for ViT-DeepSeekMoE ($k = 2, s = 1$) on ImageNet-21k→CIFAR-100 with 12 layers and 8 experts.

Table 7: Comparison of expert matching methods (Expert Weight Matching and Gate Weight Matching) across three ViT-MoE model variants evaluated on the CIFAR-10 and CIFAR-100 test sets, measuring **loss**. Metrics reported include rank and \hat{L} , defined in Section 6.3, computed across 24 permutations for 10 checkpoint pairs. All models consist of 12 layers and 4 experts.

Method	Dataset	Layer replaced	Expert Weight Matching		Gate Weight Matching	
			Rank \downarrow	$\hat{L} \downarrow$	Rank	$\hat{L} \downarrow$
MoE	CIFAR-10	1	2.50 ± 1.50	2.12 ± 0.42	3.00 ± 1.00	3.42 ± 0.55
		4	2.10 ± 0.50	1.04 ± 0.32	2.60 ± 0.92	1.54 ± 0.44
		8	2.70 ± 0.46	0.60 ± 0.27	2.90 ± 0.30	0.74 ± 0.17
		12	4.60 ± 2.00	0.13 ± 0.05	3.80 ± 1.66	0.09 ± 0.03
	CIFAR-100	1	2.80 ± 0.40	3.17 ± 0.25	2.90 ± 0.70	2.73 ± 1.03
		4	3.60 ± 1.20	1.15 ± 0.55	2.70 ± 1.00	2.03 ± 0.93
		8	3.30 ± 0.78	0.67 ± 0.14	3.20 ± 0.87	1.13 ± 0.55
		12	3.40 ± 0.92	0.07 ± 0.03	4.20 ± 0.89	0.11 ± 0.04
SMoE ($k = 2$)	CIFAR-10	1	3.00 ± 1.00	3.80 ± 1.18	3.20 ± 0.98	3.46 ± 2.94
		4	2.80 ± 0.98	1.73 ± 0.49	2.40 ± 1.56	1.64 ± 0.86
		8	2.60 ± 0.49	0.40 ± 0.46	2.60 ± 0.43	0.36 ± 0.83
		12	4.00 ± 2.45	0.06 ± 0.04	2.80 ± 1.66	0.22 ± 0.19
	CIFAR-100	1	2.80 ± 0.40	3.29 ± 3.18	2.10 ± 0.70	2.00 ± 2.02
		4	2.60 ± 1.20	1.20 ± 0.43	3.20 ± 1.47	1.91 ± 1.05
		8	3.10 ± 0.83	0.13 ± 0.07	2.70 ± 0.90	0.11 ± 0.09
		12	2.40 ± 0.92	0.07 ± 0.13	2.00 ± 0.89	0.03 ± 0.12
DeepSeekMoE ($k = 2, s = 1$)	CIFAR-10	1	3.10 ± 2.12	5.06 ± 1.22	3.60 ± 0.83	4.28 ± 2.92
		4	2.60 ± 0.77	3.32 ± 0.38	3.10 ± 1.52	2.80 ± 0.88
		8	2.30 ± 0.51	0.39 ± 0.48	2.40 ± 0.37	0.78 ± 1.00
		12	3.60 ± 3.55	0.08 ± 0.04	3.60 ± 1.48	0.23 ± 0.17
	CIFAR-100	1	4.20 ± 0.60	4.45 ± 3.40	3.70 ± 0.45	4.21 ± 2.60
		4	2.50 ± 0.78	1.34 ± 0.46	3.10 ± 1.91	2.98 ± 1.46
		8	1.90 ± 0.60	0.13 ± 0.08	3.10 ± 1.02	0.13 ± 0.16
		12	3.30 ± 1.26	0.04 ± 0.11	2.70 ± 0.89	0.02 ± 0.20

1315 G.3 Expert Matching Method

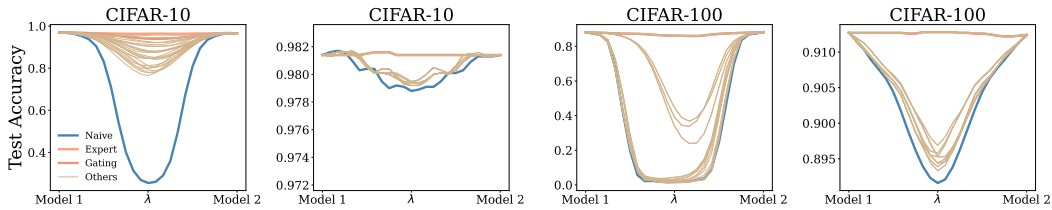


Figure 96: Accuracy curves for 12-layer ViT models incorporating a 4-expert MoE configuration in either the first layer (subplots 1 and 3) or the last layer (subplots 2 and 4), on CIFAR-10 and CIFAR-100. The curves compare two Expert Order Matching methods across 24 permutations, with Weight Matching applied post-reordering, corresponding Figure of Figure 2

Table 8: Comparison of expert matching methods (Expert Weight Matching and Gate Weight Matching) across three ViT-MoE model variants evaluated on the CIFAR-10 and CIFAR-100 test sets, measuring **accuracy**. Metrics reported include rank and \hat{L} , defined in Section 6.3, computed across 24 permutations for 10 checkpoint pairs. All models consist of 12 layers and 4 experts.

Method	Dataset	Layer replaced	Expert Weight Matching		Gate Weight Matching	
			Rank \downarrow	$\hat{L} \downarrow$	Rank	$\hat{L} \downarrow$
MoE	CIFAR-10	1	4.90 ± 2.59	1.80 ± 0.78	4.40 ± 1.80	2.00 ± 0.44
		4	3.20 ± 0.40	0.95 ± 0.97	2.90 ± 0.70	1.64 ± 0.98
		8	3.00 ± 1.80	0.03 ± 0.05	3.80 ± 1.83	0.05 ± 0.05
		12	4.00 ± 1.45	0.02 ± 0.03	2.80 ± 1.40	0.04 ± 0.06
	CIFAR-100	1	3.20 ± 1.47	2.66 ± 1.03	4.20 ± 1.00	1.51 ± 1.16
		4	2.70 ± 0.78	1.87 ± 1.17	2.90 ± 0.83	1.39 ± 1.06
		8	4.70 ± 1.42	0.08 ± 0.02	3.80 ± 1.47	0.05 ± 0.03
		12	2.40 ± 0.92	0.05 ± 0.02	3.20 ± 0.87	0.02 ± 0.02
SMoE ($k = 2$)	CIFAR-10	1	2.30 ± 2.06	3.09 ± 1.21	3.20 ± 2.75	2.24 ± 0.90
		4	2.90 ± 0.59	1.16 ± 1.00	2.80 ± 0.83	1.75 ± 1.00
		8	3.40 ± 1.47	0.07 ± 0.03	3.00 ± 1.74	0.04 ± 0.02
		12	3.60 ± 1.21	0.12 ± 0.08	4.70 ± 2.05	0.07 ± 0.03
	CIFAR-100	1	2.10 ± 1.85	3.52 ± 2.26	2.40 ± 0.93	2.05 ± 2.95
		4	2.00 ± 0.67	1.81 ± 1.28	2.90 ± 0.73	1.02 ± 2.41
		8	4.30 ± 1.60	0.09 ± 0.03	3.10 ± 1.02	0.11 ± 0.08
		12	2.60 ± 1.03	0.11 ± 0.08	2.70 ± 0.89	0.03 ± 0.06
DeepSeekMoE ($k = 2, s = 1$)	CIFAR-10	1	2.80 ± 2.20	2.95 ± 1.15	3.30 ± 2.40	2.60 ± 0.85
		4	3.10 ± 0.80	1.40 ± 1.13	2.90 ± 0.90	1.60 ± 1.30
		8	3.60 ± 1.32	0.09 ± 0.04	3.40 ± 1.10	0.07 ± 0.03
		12	3.20 ± 1.00	0.06 ± 0.02	3.00 ± 0.87	0.05 ± 0.02
	CIFAR-100	1	2.54 ± 1.90	3.20 ± 2.00	2.80 ± 2.00	2.90 ± 1.80
		4	2.90 ± 0.70	1.60 ± 1.20	3.10 ± 0.80	1.40 ± 1.00
		8	3.80 ± 1.43	0.10 ± 0.03	3.50 ± 1.20	0.08 ± 0.04
		12	2.70 ± 0.90	0.07 ± 0.02	2.90 ± 0.70	0.06 ± 0.03

Table 9: Ablation Study on varying the number of Transformer layers and all MoE variants Integration in ViT for CIFAR-10. The ratios of metrics (loss barrier, loss AUC, accuracy barrier, and accuracy AUC) compare Algorithm 1 to naive interpolation. For full table results (non-raio), kindly refer to Tables 10 and 11

Method	Number of layers	Layer replaced	Loss barrier ratio (%) ↓	Loss AUC ratio (%) ↓	Acc. barrier ratio (%) ↓	Acc. AUC ratio (%) ↓
MoE	2	1	8.54 ± 1.53	8.73 ± 1.68	8.39 ± 1.42	8.01 ± 1.39
		2	9.33 ± 2.04	8.94 ± 1.92	9.10 ± 2.19	8.83 ± 2.22
	4	1	8.29 ± 1.63	8.08 ± 1.59	7.41 ± 1.45	6.43 ± 1.09
		2	10.19 ± 3.43	10.19 ± 3.16	10.21 ± 5.26	9.08 ± 5.33
		3	8.50 ± 1.82	7.83 ± 1.89	9.82 ± 2.58	8.70 ± 2.58
		4	5.12 ± 0.67	4.60 ± 0.65	4.17 ± 1.68	4.58 ± 1.01
	6	1	4.16 ± 1.47	4.96 ± 2.31	3.64 ± 1.00	4.13 ± 1.55
		2	6.78 ± 4.80	7.71 ± 3.09	7.27 ± 5.73	8.98 ± 2.94
		3	9.69 ± 4.92	8.62 ± 5.46	9.03 ± 3.31	8.14 ± 3.57
		4	10.83 ± 5.61	12.01 ± 6.42	8.92 ± 2.21	10.00 ± 7.30
		5	9.56 ± 4.52	11.72 ± 5.59	8.72 ± 2.89	11.01 ± 3.83
		6	6.90 ± 3.54	9.67 ± 6.59	7.16 ± 3.98	9.69 ± 3.11
SMoE ($k = 2$)	2	1	8.93 ± 0.92	9.34 ± 1.23	9.36 ± 1.33	8.84 ± 1.40
		2	7.55 ± 1.39	6.62 ± 1.35	9.10 ± 2.88	7.29 ± 2.78
	4	1	8.96 ± 1.54	8.83 ± 1.64	8.66 ± 2.07	8.21 ± 1.91
		2	9.62 ± 2.68	10.07 ± 3.50	9.41 ± 3.04	9.67 ± 3.52
		3	13.31 ± 1.49	12.44 ± 2.06	11.37 ± 3.01	13.63 ± 3.17
		4	9.17 ± 1.40	9.36 ± 1.48	10.97 ± 2.77	12.52 ± 2.57
	6	1	2.37 ± 1.92	2.65 ± 1.34	2.25 ± 2.71	4.67 ± 2.38
		2	10.09 ± 3.11	9.26 ± 3.19	8.35 ± 3.37	8.36 ± 2.38
		3	11.22 ± 6.25	10.17 ± 8.52	12.10 ± 7.04	12.80 ± 10.77
		4	8.76 ± 2.06	12.16 ± 2.22	13.09 ± 5.49	10.85 ± 7.21
		5	14.39 ± 7.82	14.66 ± 9.10	11.82 ± 9.19	12.45 ± 7.85
		6	6.38 ± 5.01	9.56 ± 3.79	11.96 ± 11.93	10.93 ± 15.06
DeepSeekMoE ($k = 2, r = 1$)	2	1	8.94 ± 1.95	9.12 ± 2.10	7.96 ± 1.59	7.88 ± 1.90
		2	9.25 ± 0.37	8.27 ± 0.44	12.00 ± 1.35	11.14 ± 1.73
	4	1	8.65 ± 1.01	8.56 ± 1.19	7.17 ± 0.99	6.46 ± 0.88
		2	10.80 ± 3.77	10.88 ± 3.69	10.45 ± 4.91	9.92 ± 4.83
		3	7.90 ± 1.79	8.82 ± 1.87	10.21 ± 2.44	9.02 ± 2.52
		4	8.35 ± 1.20	8.29 ± 1.39	9.29 ± 2.75	7.04 ± 3.14
	6	1	6.29 ± 5.00	6.17 ± 3.48	8.07 ± 5.34	8.35 ± 5.65
		2	8.03 ± 1.79	9.16 ± 2.10	7.88 ± 3.66	8.54 ± 3.93
		3	9.27 ± 4.25	9.29 ± 3.19	12.93 ± 2.18	11.56 ± 5.58
		4	10.80 ± 2.18	12.72 ± 1.52	7.34 ± 3.64	8.99 ± 5.10
		5	13.77 ± 2.27	13.03 ± 7.01	11.14 ± 5.47	13.28 ± 3.32
		6	9.82 ± 2.99	8.38 ± 4.81	8.92 ± 4.22	8.14 ± 3.48

Table 10: Ablation study results on CIFAR-10 for LMC in ViT with layer replacements of all MoE variants, reporting test set **loss**. We evaluate the loss metrics for linear interpolation using our proposed Algorithm 1 against naive linear interpolation. Metrics include the loss barrier and loss Area Under the Curve (AUC), with AUC computed relative to the straight line connecting the two model endpoints. Each experiment is repeated *five* times, and LMC is performed on *ten* model pairs. Additionally, we report loss values across five models to provide enhanced context. To improve readability, all metrics are scaled up by a factor of 10^2 .

Method	Number of layers	Layer replaced	WM loss barrier ↓	Naive loss barrier ↓	WM loss AUC ↓	Naive loss AUC ↓	Loss value ↓
MoE	2	1	3.14 ± 0.55	34.22 ± 1.25	1.45 ± 0.28	16.11 ± 0.74	99.04 ± 0.92
		2	2.84 ± 0.44	33.54 ± 1.12	1.38 ± 0.15	17.53 ± 0.57	104.52 ± 0.81
	4	1	9.48 ± 1.36	115.56 ± 10.19	4.25 ± 0.68	53.07 ± 4.54	97.00 ± 0.59
		2	6.59 ± 2.12	64.73 ± 2.05	3.71 ± 1.11	36.39 ± 0.96	91.74 ± 0.87
		3	4.20 ± 0.90	48.23 ± 0.92	2.21 ± 0.53	28.23 ± 0.50	88.06 ± 1.30
		4	2.76 ± 0.38	53.92 ± 0.80	1.42 ± 0.22	31.01 ± 0.54	86.15 ± 0.33
	6	1	1.42 ± 0.44	37.47 ± 8.47	0.74 ± 0.30	18.73 ± 3.30	90.35 ± 1.15
		2	0.47 ± 0.32	7.09 ± 0.60	0.27 ± 0.10	3.53 ± 0.30	88.73 ± 0.39
		3	0.33 ± 0.17	3.42 ± 0.22	0.13 ± 0.08	1.60 ± 0.14	88.26 ± 0.50
		4	0.14 ± 0.08	1.36 ± 0.33	0.07 ± 0.04	0.64 ± 0.16	86.48 ± 0.17
		5	0.14 ± 0.06	1.68 ± 0.44	0.08 ± 0.03	0.83 ± 0.21	85.92 ± 0.46
		6	0.23 ± 0.11	2.46 ± 0.29	0.12 ± 0.04	1.26 ± 0.21	88.61 ± 0.33
SMoE ($k = 2$)	2	1	2.95 ± 0.24	33.17 ± 1.01	1.39 ± 0.16	14.95 ± 0.51	98.88 ± 0.86
		2	2.52 ± 0.46	33.46 ± 0.67	1.20 ± 0.25	18.12 ± 0.48	105.59 ± 0.37
	4	1	10.45 ± 1.99	116.54 ± 8.34	4.75 ± 1.01	53.65 ± 3.97	96.73 ± 0.85
		2	8.31 ± 0.83	86.88 ± 3.25	4.73 ± 0.57	47.40 ± 1.91	90.80 ± 1.02
		3	4.73 ± 0.62	35.49 ± 1.34	2.59 ± 0.49	20.78 ± 0.80	86.29 ± 0.24
		4	4.48 ± 0.70	48.84 ± 0.68	2.68 ± 0.44	28.62 ± 0.44	82.33 ± 0.37
	6	1	2.12 ± 0.67	34.23 ± 7.22	1.28 ± 0.35	12.59 ± 2.67	90.25 ± 1.00
		2	0.75 ± 0.29	7.28 ± 0.83	0.34 ± 0.15	3.59 ± 0.43	88.99 ± 0.65
		3	0.34 ± 0.19	3.28 ± 0.30	0.15 ± 0.06	1.76 ± 0.13	87.69 ± 0.52
		4	0.26 ± 0.14	2.78 ± 0.38	0.14 ± 0.07	1.35 ± 0.19	86.08 ± 0.40
		5	0.22 ± 0.11	2.50 ± 0.04	0.11 ± 0.05	1.25 ± 0.02	86.08 ± 0.34
		6	0.23 ± 0.18	3.63 ± 0.17	0.13 ± 0.08	1.93 ± 0.10	88.38 ± 0.36
DeepSeekMoE ($k = 2, s = 1$)	2	1	3.00 ± 0.61	33.63 ± 1.30	1.37 ± 0.27	15.14 ± 0.73	98.41 ± 1.06
		2	3.07 ± 0.21	33.16 ± 1.26	1.50 ± 0.11	18.18 ± 0.72	105.92 ± 0.53
	4	1	10.57 ± 1.57	122.13 ± 10.50	4.91 ± 0.84	57.32 ± 4.89	97.71 ± 1.26
		2	5.85 ± 2.14	54.13 ± 2.27	2.90 ± 1.26	30.17 ± 1.35	90.89 ± 0.52
		3	2.67 ± 0.65	33.36 ± 1.47	1.35 ± 0.40	19.48 ± 0.89	86.37 ± 0.43
		4	4.11 ± 0.61	49.18 ± 0.95	2.37 ± 0.41	28.49 ± 0.55	82.26 ± 0.39
	6	1	2.25 ± 0.82	39.42 ± 6.23	1.14 ± 0.38	18.88 ± 3.32	91.32 ± 0.25
		2	2.88 ± 0.57	28.53 ± 1.26	1.56 ± 0.32	15.97 ± 0.69	89.50 ± 1.05
		3	0.56 ± 0.12	7.90 ± 0.18	0.26 ± 0.06	3.90 ± 0.12	88.22 ± 0.57
		4	0.24 ± 0.16	3.07 ± 0.34	0.12 ± 0.03	1.53 ± 0.14	86.09 ± 0.40
		5	0.32 ± 0.17	3.39 ± 0.31	0.16 ± 0.08	1.63 ± 0.13	86.10 ± 0.51
		6	0.23 ± 0.08	2.61 ± 0.53	0.14 ± 0.04	1.43 ± 0.32	89.27 ± 0.57

Table 11: Ablation study results on CIFAR-10 for LMC in ViT with layer replacements of all MoE variants, reporting test set **accuracy**. We evaluate the accuracy metrics for linear interpolation using our proposed Algorithm 1 against naive linear interpolation. Metrics include the accuracy barrier and accuracy Area Under the Curve (AUC), with AUC computed relative to the straight line connecting the two model endpoints. Each experiment is repeated *five* times, and LMC is performed on *ten* model pairs. Additionally, we report accuracy values across five models to provide enhanced context. To improve readability, all metrics are scaled up by a factor of 10^2 .

Method	Number of layers	Layer replaced	WM acc. barrier ↓	Naive acc. barrier ↓	WM acc. AUC ↓	Naive acc. barrier ↓	Acc. value ↑
MoE	2	1	0.91 ± 0.28	12.33 ± 1.28	0.44 ± 0.15	5.89 ± 0.54	66.53 ± 3.30
		2	0.76 ± 0.21	7.85 ± 2.31	0.33 ± 0.08	3.77 ± 1.11	63.65 ± 0.47
	4	1	2.10 ± 0.42	28.36 ± 1.85	0.86 ± 0.16	13.43 ± 0.99	72.48 ± 0.68
		2	1.41 ± 0.71	13.83 ± 0.53	0.75 ± 0.41	7.88 ± 0.29	73.43 ± 0.61
		3	0.67 ± 0.17	6.84 ± 0.24	0.35 ± 0.09	4.00 ± 0.12	73.64 ± 0.20
		4	0.23 ± 0.09	5.48 ± 0.19	0.14 ± 0.03	3.05 ± 0.14	73.58 ± 0.19
	6	1	0.50 ± 0.13	13.95 ± 2.20	0.24 ± 0.07	5.71 ± 0.93	74.29 ± 0.35
		2	0.35 ± 0.10	2.49 ± 0.36	0.12 ± 0.04	1.23 ± 0.21	74.65 ± 0.44
		3	0.14 ± 0.05	1.48 ± 0.28	0.05 ± 0.02	0.72 ± 0.14	74.72 ± 0.46
		4	0.05 ± 0.08	1.33 ± 0.43	0.03 ± 0.05	0.68 ± 0.24	74.46 ± 0.12
		5	0.05 ± 0.03	0.59 ± 0.05	0.02 ± 0.01	0.28 ± 0.02	74.62 ± 0.08
		6	0.06 ± 0.04	0.61 ± 0.12	0.03 ± 0.02	0.30 ± 0.05	74.25 ± 0.21
SMoE ($k = 2$)	2	1	1.07 ± 0.15	11.42 ± 0.22	0.46 ± 0.07	5.20 ± 0.09	65.27 ± 0.69
		2	0.58 ± 0.17	6.41 ± 0.24	0.26 ± 0.10	3.63 ± 0.13	62.58 ± 0.16
	4	1	2.48 ± 0.66	28.55 ± 1.55	1.12 ± 0.29	13.52 ± 0.81	72.26 ± 0.60
		2	1.78 ± 0.23	19.00 ± 0.71	0.97 ± 0.15	10.16 ± 0.36	73.12 ± 0.40
		3	0.87 ± 0.17	4.98 ± 0.32	0.38 ± 0.10	2.78 ± 0.22	73.58 ± 0.33
		4	0.79 ± 0.14	5.29 ± 0.31	0.40 ± 0.09	3.15 ± 0.15	73.60 ± 0.44
	6	1	0.30 ± 0.38	12.61 ± 2.27	0.28 ± 0.11	5.11 ± 0.86	74.09 ± 0.28
		2	0.29 ± 0.16	2.85 ± 0.37	0.13 ± 0.08	1.48 ± 0.22	74.81 ± 0.47
		3	0.18 ± 0.14	1.36 ± 0.18	0.10 ± 0.06	0.95 ± 0.10	74.97 ± 0.30
		4	0.18 ± 0.07	1.35 ± 0.05	0.08 ± 0.02	0.76 ± 0.03	74.60 ± 0.24
		5	0.19 ± 0.09	1.45 ± 0.08	0.09 ± 0.04	0.77 ± 0.05	74.72 ± 0.24
		6	0.19 ± 0.08	1.60 ± 0.14	0.08 ± 0.02	0.89 ± 0.09	74.15 ± 0.28
DeepSeekMoE ($k = 2, s = 1$)	2	1	0.92 ± 0.18	11.64 ± 0.48	0.42 ± 0.10	5.29 ± 0.25	65.43 ± 0.30
		2	0.78 ± 0.11	6.45 ± 0.31	0.41 ± 0.08	3.67 ± 0.21	62.55 ± 0.27
	4	1	2.07 ± 0.31	28.85 ± 1.64	0.90 ± 0.13	13.93 ± 0.91	71.57 ± 0.84
		2	1.77 ± 0.56	12.14 ± 0.50	0.67 ± 0.32	6.74 ± 0.24	73.70 ± 0.14
		3	1.01 ± 0.14	9.80 ± 0.28	0.55 ± 0.08	4.91 ± 0.17	73.72 ± 0.20
		4	0.52 ± 0.15	5.60 ± 0.14	0.23 ± 0.10	3.22 ± 0.11	73.89 ± 0.30
	6	1	2.43 ± 0.73	30.80 ± 2.43	0.99 ± 0.31	14.06 ± 0.90	74.26 ± 0.19
		2	0.95 ± 0.19	12.63 ± 0.43	0.47 ± 0.10	7.37 ± 0.18	74.51 ± 0.14
		3	0.39 ± 0.13	4.85 ± 0.14	0.20 ± 0.06	2.44 ± 0.09	74.98 ± 0.19
		4	0.14 ± 0.07	2.29 ± 0.08	0.10 ± 0.04	1.13 ± 0.03	74.70 ± 0.16
		5	0.12 ± 0.03	1.39 ± 0.08	0.08 ± 0.02	0.66 ± 0.04	74.62 ± 0.14
		6	0.16 ± 0.07	1.74 ± 0.12	0.06 ± 0.04	0.88 ± 0.05	74.31 ± 0.49

1317 **H Broader Impact**

1318 The investigation into Linear Mode Connectivity (LMC) within Mixture-of-Experts (MoE) architec-
1319 tures presents both potential positive and negative societal implications. On the positive side, this
1320 work could lead to more efficient training methods for MoE models, enhancing their scalability and
1321 computational efficiency. Such advancements may democratize access to advanced AI technologies,
1322 enabling researchers and developers with limited resources to leverage powerful models for innova-
1323 tive applications. Furthermore, the insights gained into the optimization dynamics and functional
1324 landscape of neural networks could contribute to the development of models that generalize more
1325 effectively, thereby improving the reliability and robustness of AI systems in critical domains such
1326 as healthcare, finance, and autonomous systems. Additionally, this research enriches the broader
1327 AI community’s understanding of neural network loss landscapes, fostering further theoretical and
1328 practical advancements.

1329 It is important to note that while these potential impacts are significant, the primary contribution of
1330 this work lies in its theoretical and foundational insights. The actual societal ramifications will largely
1331 depend on how these insights are applied and governed in practical settings. As such, responsible
1332 development and deployment practices, coupled with ongoing research into bias mitigation and ethical
1333 AI, are essential to maximizing the benefits and minimizing the risks associated with advancements
1334 in MoE architectures.

1335 NeurIPS Paper Checklist

1336 1. Claims

1337 Question: Do the main claims made in the abstract and introduction accurately reflect the
1338 paper’s contributions and scope?

1339 Answer: [\[Yes\]](#)

1340 Justification: The claims made in the abstract and introduction are clearly stated in the
1341 **Contribution** in the Introduction. We provide mathematical contexts in Section 2, which
1342 provides background on LMC and MoE architectures. We introduce the concept of the
1343 weight space of MoE architectures and define a group action on this space that preserves the
1344 functional behavior of MoE models, aligning with the claim of defining such a space and
1345 action in Section 3. We present two core results concerning functional equivalence in MoE
1346 models, demonstrating that the proposed group action characterizes all inherent symmetries
1347 of the MoE gating mechanism in Section 4. We develop a Weight Matching algorithm that
1348 enables alignment between independently trained MoEs in Section 5. We provide empirical
1349 evidence of LMC across a wide range of MoE configurations and additional experiments to
1350 support our work in Section 6.

1351 Guidelines:

- 1352 • The answer NA means that the abstract and introduction do not include the claims
1353 made in the paper.
- 1354 • The abstract and/or introduction should clearly state the claims made, including the
1355 contributions made in the paper and important assumptions and limitations. A No or
1356 NA answer to this question will not be perceived well by the reviewers.
- 1357 • The claims made should match theoretical and experimental results, and reflect how
1358 much the results can be expected to generalize to other settings.
- 1359 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
1360 are not attained by the paper.

1361 2. Limitations

1362 Question: Does the paper discuss the limitations of the work performed by the authors?

1363 Answer: [\[Yes\]](#)

1364 Justification: The limitations are discussed in the the Conclusion.

1365 Guidelines:

- 1366 • The answer NA means that the paper has no limitation while the answer No means that
1367 the paper has limitations, but those are not discussed in the paper.
- 1368 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 1369 • The paper should point out any strong assumptions and how robust the results are to
1370 violations of these assumptions (e.g., independence assumptions, noiseless settings,
1371 model well-specification, asymptotic approximations only holding locally). The authors
1372 should reflect on how these assumptions might be violated in practice and what the
1373 implications would be.
- 1374 • The authors should reflect on the scope of the claims made, e.g., if the approach was
1375 only tested on a few datasets or with a few runs. In general, empirical results often
1376 depend on implicit assumptions, which should be articulated.
- 1377 • The authors should reflect on the factors that influence the performance of the approach.
1378 For example, a facial recognition algorithm may perform poorly when image resolution
1379 is low or images are taken in low lighting. Or a speech-to-text system might not be
1380 used reliably to provide closed captions for online lectures because it fails to handle
1381 technical jargon.
- 1382 • The authors should discuss the computational efficiency of the proposed algorithms
1383 and how they scale with dataset size.
- 1384 • If applicable, the authors should discuss possible limitations of their approach to
1385 address problems of privacy and fairness.

- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [\[Yes\]](#)

Justification: All theoretical results in the paper are given together with the full set of assumptions and complete/correct proofs (See Sections 3, 4, 5 and Appendices A, B, C, D in our manuscript).

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [\[Yes\]](#)

Justification: We provide the experiment details in the Experiment Details Section (Appendix F) in the Appendix of our manuscript. We also provide the source code so that the results in the paper can be easily reproduced.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.

- (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
- (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
- (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: We provide the source code in the Supplemental Materials so that the results in the paper can be easily reproduced. We verify our proposed methods using public benchmarks (See the Experimental Results Section, i.e., Section 6, in our manuscript)

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: We specify all the training and test details necessary to understand the results in the Experimental Results Section (Section 6) in the maintext and the Experiment Details Section (Appendix F) in the Appendix of our manuscript.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: We report error bars suitably and correctly defined of the experiments.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: We provide sufficient information on the computer resources for all experiments in our Experimental Results Section (Section 6) and Appendix F.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

Answer: [Yes]

Justification: The research conducted in the paper conforms, in every respect, with the NeurIPS Code of Ethics.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.

- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [\[Yes\]](#)

Justification: We discuss broader impacts in Appendix H.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [\[NA\]](#)

Justification: The paper poses no such risks.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [\[Yes\]](#)

Justification: We cite the githubs we use and the baselines we compare with in our manuscript. All the assets used in the paper are properly cited.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. New assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [\[Yes\]](#)

Justification: We include details about training and implementation in Appendix F, and code in supplementary materials. The new assets provided in the paper are well documented, and the documentation is provided alongside the assets.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. Crowdsourcing and research with human subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [\[NA\]](#)

Justification: The paper does not involve crowdsourcing or research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. Institutional review board (IRB) approvals or equivalent for research with human subjects

1647 Question: Does the paper describe potential risks incurred by study participants, whether
 1648 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
 1649 approvals (or an equivalent approval/review based on the requirements of your country or
 1650 institution) were obtained?

1651 Answer: [NA]

1652 Justification: The paper does not involve research with human subjects.

1653 Guidelines:

- 1654 • The answer NA means that the paper does not involve crowdsourcing nor research with
- 1655 human subjects.
- 1656 • Depending on the country in which research is conducted, IRB approval (or equivalent)
- 1657 may be required for any human subjects research. If you obtained IRB approval, you
- 1658 should clearly state this in the paper.
- 1659 • We recognize that the procedures for this may vary significantly between institutions
- 1660 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
- 1661 guidelines for their institution.
- 1662 • For initial submissions, do not include any information that would break anonymity (if
- 1663 applicable), such as the institution conducting the review.

1664 **16. Declaration of LLM usage**

1665 Question: Does the paper describe the usage of LLMs if it is an important, original, or

1666 non-standard component of the core methods in this research? Note that if the LLM is used

1667 only for writing, editing, or formatting purposes and does not impact the core methodology,

1668 scientific rigor, or originality of the research, declaration is not required.

1669 Answer: [NA]

1670 Justification: The core methodological contributions of this research does NOT rely on

1671 LLMs in any any important, original, or non-standard way.

1672 Guidelines:

- 1673 • The answer NA means that the core method development in this research does not
- 1674 involve LLMs as any important, original, or non-standard components.
- 1675 • Please refer to our LLM policy (<https://neurips.cc/Conferences/2025/LLM>)
- 1676 for what should or should not be described.