

Embedding Principle: A Hierarchical Structure of Loss Landscape of Deep Neural Networks

Yaoyu Zhang ^{* 1}, Yuqing Li ^{† 2}, Zhongwang Zhang ^{‡ 3}, Tao Luo ^{§ 4}, and
Zhi-Qin John Xu ^{¶ 5}

¹School of Mathematical Sciences, Institute of Natural Sciences, MOE-LSC and Qing Yuan Research Institute, Shanghai Jiao Tong University, Shanghai Center for Brain Science and Brain-Inspired Technology, Shanghai, 200240, China.

²School of Mathematical Sciences, CMA-Shanghai, Shanghai Jiao Tong University, Shanghai, 200240, China.

³School of Mathematical Sciences, Institute of Natural Sciences, Shanghai Jiao Tong University, Shanghai, 200240, China.

⁴School of Mathematical Sciences, CMA-Shanghai, Institute of Natural Sciences, MOE-LSC and Qing Yuan Research Institute, Shanghai Jiao Tong University, Shanghai, 200240, China.

⁵Institute of Natural Sciences, School of Mathematical Sciences, MOE-LSC and Qing Yuan Research Institute, Shanghai Jiao Tong University, Shanghai, 200240, China.

Abstract. We prove a general Embedding Principle of loss landscape of deep neural networks (NNs) that unravels a hierarchical structure of the loss landscape of NNs, i.e., loss landscape of an NN *contains* all critical points of all the narrower NNs. This result is obtained by constructing a class of critical embeddings which map any critical point of a narrower NN to a critical point of the target NN with the same output function. By discovering a wide class of general compatible critical embeddings, we provide a gross estimate of the dimension of critical submanifolds embedded from critical points of narrower NNs. We further prove an irreversibility property of any critical embedding that the number of negative/zero/positive eigenvalues of the Hessian matrix of a critical point may increase but never decrease as an NN becomes wider through the embedding. Using a special realization of general compatible critical embedding, we prove a stringent necessary condition for being a “truly-bad” critical point that never becomes a strict-saddle point through any critical embedding. This result implies the commonplace of strict-saddle points in wide NNs, which may be an important reason underlying the easy optimization of wide NNs widely observed in practice.

Keywords:

Neural network,
Loss landscape,
Critical point,
Embedding principle.

Article Info.:

Volume: 1
Number: 1
Pages: 60- 113
Date: March/2022
doi.org/10.4208/jml.220108

Article History:

Received: 2/12/2021
Accepted: 19/3/2022

Communicated by:

Arnulf Jentzen

1 Introduction

The loss landscape of a deep neural network (NN) is important to both its optimization dynamics and generalization performance, hence is a key issue in deep learning theory. It

*Corresponding author. zhyy.sjtu@sjtu.edu.cn.

†liyqing_551@sjtu.edu.cn.

‡0123zzw666@sjtu.edu.cn.

§luotao41@sjtu.edu.cn.

¶Corresponding author. xuzhiqin@sjtu.edu.cn.

has been realized for a long time that it is important to quantify exactly how the loss landscape looks like [1,2]. This problem is difficult since various visualization methods show that the NN loss landscape is very complicated [3,4]. Moreover, its non-convexity, high dimensionality and the dependence on data, model and the specific form of loss function make it very difficult to obtain a general understanding through empirical study. Therefore, though it has been extensively studied over the years, it remains an open problem to provide a clear picture about the general structure of a DNN loss landscape, e.g., critical points/submanifolds, their output functions and other properties.

Our work is inspired by the following empirical observations. From the aspect of optimization, it is often observed that wider NNs are easier for training. This phenomenon not only holds in a neural tangent kernel (NTK) regime [5], where the gradient descent training can find the global minimum with a linear convergence rate [6–9], but also happens in highly nonlinear regimes beyond NTK [10,11]. From the aspect of generalization, the puzzle that over-parameterized NNs often generalize well seems to contradict the conventional learning theory [12,13]. The frequency principle [14–19] shows that NNs, over-parameterized or not, tend to fit the training data by a low-frequency function, which suggests that the learned function by an NN is often of much lower complexity than the NN’s capacity. Specifically, with small initialization, e.g., in a condensed regime, weights of an NN are empirically found to condense on isolated directions resulting in an output function mimicking that of a narrower NN [11,20]. These observations raise a question that in which sense learning of a wide NN is not drastically different from a narrower NN despite potentially huge difference in their numbers of parameters. From the aspect of pruning, empirical works propose a “lottery ticket hypothesis” that a substantially smaller sub-network can achieve the same accuracy as the original large network [21]. However, it is not yet clear about the mechanism of redundancy in a learned wide NN, which makes a drastic pruning possible in practice.

All above empirical observations, though relevant to different aspects of NN, are in essence pointing towards an intrinsic similarity between narrow and wide NNs. In this work, focusing on the loss landscape, we address the following problem: What are the relations of critical points and the corresponding output functions of loss landscape among NNs with different widths. The significance of studying the critical points and their output functions of the NN loss landscape is as follows. From the optimization perspective, NN parameters trained by gradient descent provably converge to a critical point, which in general is not necessarily a global minimum or local minimum. Moreover, even for these saddle points which can be escaped, e.g., strict-saddle points [22], they may still attract the training trajectory (points nearby with a higher loss will first come close and then move away), contributing to the implicit regularization of NNs, say towards a simpler fitting (i.e., a fitting that can be realized by a narrower NN). We are specifically interested in output functions corresponding to critical points, named as *critical functions* for convenience. Studying these critical functions that potentially attract the learning of NN is clearly important for a deeper understanding of the learning process of an NN.

Our key finding in this work is the following general principle about critical points/functions of NN loss landscape intuitively stated as follows:

Embedding principle: the loss landscape of an NN contains all critical points/functions of all the narrower NNs.

The Embedding Principle shows that any NN loss landscape contains a hierarchical structure of critical points/functions with different complexities from NNs of different widths. Specifically, it ensures existence of “simple” critical functions that can be represented by narrow NNs. Therefore, combining with the phenomenon of Frequency Principle and condensation, we conjecture that nonlinear training of NNs may be implicitly biased towards these “simple” critical functions. We will carefully look into this conjecture in our future works.

To prove the Embedding Principle, we first construct one-step critical embeddings which map any parameters of a narrow network to that of an one-neuron wider NN preserving the output function and criticality. With these embeddings, critical points of a narrow network loss landscape is mapped to 1-d critical affine subspaces of an one-neuron wider network loss landscape with the same output function. These one-step critical embeddings are constructed by adding a null neuron or splitting an existing neuron. By composition of one-step embeddings, any critical point of a narrow network loss landscape can be mapped to a critical point of any wider network loss landscape preserving the output function. Importantly, we further propose a wide class of general compatible critical embeddings, where one-step embeddings or their composition are its special cases. Note that, all critical points of a wide NN embedded from a critical point of a narrower NN by all possible general compatible critical embeddings, form high-dimensional critical submanifolds which in general are not affine subspaces for three-layer or deeper NNs.

The critical embeddings naturally link critical points of NNs of different widths, thus providing a means to track how properties of these critical points may change as the width of the NN increases. Using these critical embeddings as a tool, we obtain rich information about the general structure of an NN loss landscape.

We show that the degeneracy of a critical point substantially increases when it is embedded to a wider network, due to the fact that a critical point can be mapped to a high-dimensional critical submanifold through a class of critical embeddings. This degeneracy of critical points arises from the neuron redundancy of the wide NN in representing certain simple critical functions from narrower NNs, which is different from over-parameterization induced degeneracy studied in [23]. We also study the property of Hessian of critical points through critical embedding, e.g., the number of its negative eigenvalues, which determines whether the corresponding critical point is a strict-saddle that enables easy optimization [22]. We prove an irreversibility of critical embedding that the number of negative eigenvalues of Hessian matrix may increase but never decrease as an NN becomes wider through critical embedding. Moreover, we introduce a notion of “truly-bad” critical point which never becomes a strict-saddle point through any critical embedding. We prove a stringent necessary condition for being a “truly-bad” critical point that requires an important ingredient of its Hessian matrix being a zero matrix. This result implies the commonplace of strict-saddle points in the high-dimensional critical submanifolds of wide NNs, which may be an important reason underlying the easy optimization of wide NNs widely observed in practice.

In summary, the following general understanding of an NN loss landscape is obtained

by the embedding principle in this work:

- (i) It contains a hierarchical structure of critical points/functions with different complexities from that of all narrower NNs;
- (ii) Critical functions from narrower NNs in general forms a high-dimensional critical submanifold with a gross estimate of the dimension: $K + \sum_{k \in [L]} K_l K_{l-1}$, where K_l is the difference in neuron number in layer l between the target NN and the narrower NN.
- (iii) If it has critical points other than the global minima that are not strict-saddle points, they mostly can become strict-saddle points in wider NNs through critical embedding, which means the embedded critical points become more optimization-friendly in a wider NN. Remark that, other than the embedded critical points, further study is needed to better understand the optimization property of these new critical points arising in wider NNs.

2 Related works

In our previous conference paper [24], we prove the Embedding Principle inspired by experimental observations and study one-step embeddings and their multi-step composition for general deep NNs. Note that, similar results on composition embedding are studied in other works, e.g., for shallow NNs [25] and deep NNs [26, 27]. As a comprehensive extension of [24], this work mathematically formalizes the notion and the study of critical embedding, proposing the general compatible embedding for the first time, and further analyzing the transition of hessian of critical points through the critical embedding. Other than above works focusing on universal properties of the NN loss landscape, many researches study the loss landscape in detail for specific settings, e.g., shallow NNs with specific activations [28–30].

Simple gradient-descent-based optimization on the complex loss landscape of NN [1, 4] often finds solutions that generalize well. Many works study the geometry of the NN loss landscape at critical points in relation to its generalization ability. For example, empirical works show that SGD [31] and dropout [32] training can find a flat minimizer, which may explain why such stochastic training can find solution that generalize better. [33] further suggest that the volume of basin of attraction of good (flat) minima may dominate over that of poor (sharp) minima in practical problems. [34] show that at a local minimum there exist many asymmetric directions such that the loss increases abruptly along one side, and slowly along the opposite side. [35] prove that for any multi-layer network with generic input data and non-linear activation functions, sub-optimal local minima can exist, no matter how wide the network is. When the network width increases towards infinity, the loss landscape may become simpler and the training can avoid spurious valleys with high probability in an over-parameterized regime [36]. In an extremely over-parameterized regime with a large initialization, i.e., the linear regime identified in [11] with the NTK regime as its special case [5], the gradient descent training can find the global minimum with a linear convergence rate [6–9, 11].

The starting point of this work originates from our work in [11], where we identify a highly nonlinear condensed regime far beyond the NTK regime that weights condense in isolated directions during the training. Moreover, neural networks of different width often exhibit similar condensed behavior, e.g., stagnating at similar loss with almost the same output function, which is illustrated in experiments in our conference paper [24]. The condensation is a highly nonlinear feature learning process important to implicit regularization and generalization of NNs. The condensation transforms a large network to a network of only a few effective neurons, leading to an output function with low complexity. Such learning process is consistent with another line of research, that is, the complexity of NN output gradually increases during the training [14–18, 37–42]. For example, the Frequency Principle [14, 15] states that NNs often fit target functions from low to high frequencies during the training. A series of works study the mechanism of condensation at an initial training stage, such as for ReLU network [20, 43] and network with continuously differentiable activation functions [44].

This work in some sense serves as our attempt to uncover the theoretical structure underlying the condensation phenomenon from the perspective of loss function by proving a general Embedding Principle. In another aspect, the condensation phenomenon also confirms the value of Embedding Principle in understanding the highly nonlinear training behavior in practice.

The Embedding Principle provides a structural mechanism underlying the degeneracy as a very common property for critical points [45, 46]. Thus it complements the understanding that global minima of NNs typically form a high dimensional manifold due to over-parameterization [23].

3 Preliminary

3.1 Deep neural networks

Consider L -layer ($L \geq 2$) fully-connected NNs with a general differentiable activation function. We regard the input as the 0-th layer and the output as the L -th layer. Let m_l be the number of neurons in the l -th layer. In particular, we also set $m_0 = d$ and $m_L = d'$. For any $i, k \in \mathbb{N}$ and $i < k$, we denote $[i : k] = \{i, i + 1, \dots, k\}$. In particular, we denote $[k] := \{1, 2, \dots, k\}$. For a matrix \mathbf{A} , we use $(\mathbf{A})_{i,j}$ to denote its (i, j) -th entry. We will also define $(\mathbf{A})_{i,[j:k]} := ((\mathbf{A})_{i,j}, (\mathbf{A})_{i,j+1}, \dots, (\mathbf{A})_{i,k})$ as part of the i -th row vector. Similarly, $(\mathbf{A})_{[j:k],i}$ is a part of the i -th column vector. For a vector \mathbf{a} , we use $(\mathbf{a})_i$ to denote its i -th entry, we also define $(\mathbf{a})_{[j:k]} := ((\mathbf{a})_j, (\mathbf{a})_{j+1}, \dots, (\mathbf{a})_k)$ as part of the vector. Given weights $\mathbf{W}^{[l]} \in \mathbb{R}^{m_l \times m_{l-1}}$ and bias $\mathbf{b}^{[l]} \in \mathbb{R}^{m_l}$ for $l \in [L]$, we define the collection of parameters $\boldsymbol{\theta}$ as a $2L$ -tuple (an ordered list of $2L$ elements) whose elements are matrices or vectors

$$\boldsymbol{\theta} := (\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_L) := (\mathbf{W}^{[1]}, \mathbf{b}^{[1]}, \dots, \mathbf{W}^{[L]}, \mathbf{b}^{[L]}), \quad (3.1)$$

where the l -th layer parameters of $\boldsymbol{\theta}$ is the ordered pair $\boldsymbol{\theta}_l := (\mathbf{W}^{[l]}, \mathbf{b}^{[l]})$, $l \in [L]$. We may misuse our notations and do not distinguish $\boldsymbol{\theta}$ from its vectorization $\text{vec}(\boldsymbol{\theta}) \in \mathbb{R}^M$ with

$M := \sum_{l=0}^{L-1} (m_l + 1)m_{l+1}$. Moreover, we call the collection of tuples of length $2L$ **the tuple class**, whose elements are matrices $\{\mathbf{W}^{[l]}\}_{l=1}^L$ with $\mathbf{W}^{[l]} \in \mathbb{R}^{m_l \times m_{l-1}}$, or vectors $\{\mathbf{b}^{[l]}\}_{l=1}^L$ with $\mathbf{b}^{[l]} \in \mathbb{R}^{m_l}$, and denoted by $\text{Tuple}_{\{m_0, \dots, m_L\}}$, i.e.,

$$\begin{aligned} & \text{Tuple}_{\{m_0, \dots, m_L\}} \\ & := \left\{ \boldsymbol{\theta} \mid \boldsymbol{\theta} = \left(\mathbf{W}^{[1]}, \mathbf{b}^{[1]}, \dots, \mathbf{W}^{[L]}, \mathbf{b}^{[L]} \right), \mathbf{W}^{[l]} \in \mathbb{R}^{m_l \times m_{l-1}}, \mathbf{b}^{[l]} \in \mathbb{R}^{m_l}, l \in [L] \right\}. \end{aligned}$$

Since the tuple class inherits the structure of Euclidean spaces, obviously it is a linear space. We set $\mathbf{0}$ as the zero element in the tuple class, i.e.

$$\mathbf{0} = \left(\mathbf{0}_{m_1 \times m_0}, \mathbf{0}_{m_1 \times 1}, \dots, \mathbf{0}_{m_L \times m_{L-1}}, \mathbf{0}_{m_L \times 1} \right) \in \text{Tuple}_{\{m_0, \dots, m_L\}},$$

and we abuse the notation $\mathbf{0}$ from time to time to denote zero elements belonging to different tuple classes.

We further define the upper bracket $[L-1]$ by limiting ourselves to the first $2L-2$ element of the tuple, i.e.,

$$\boldsymbol{\theta}^{[L-1]} := \left(\mathbf{W}^{[1]}, \mathbf{b}^{[1]}, \dots, \mathbf{W}^{[L-2]}, \mathbf{b}^{[L-2]}, \mathbf{W}^{[L-1]}, \mathbf{b}^{[L-1]} \right) \in \text{Tuple}_{\{m_0, \dots, m_{L-2}, m_{L-1}\}}, \quad (3.2)$$

and given $M^{[L-1]} := \sum_{l=0}^{L-2} (m_l + 1)m_{l+1}$, $\boldsymbol{\theta}^{[L-1]} \in \mathbb{R}^{M^{[L-1]}}$.

Given parameters $\boldsymbol{\theta}$, the neural network function $f_{\boldsymbol{\theta}}(\cdot)$ can be defined in a recursive way. First, we write $f_{\boldsymbol{\theta}}^{[0]}(\mathbf{x}) := \mathbf{x}$ for the input $\mathbf{x} \in \mathbb{R}^d$, then for $l \in [L-1]$, $f_{\boldsymbol{\theta}}^{[l]}$ is defined recursively as

$$f_{\boldsymbol{\theta}}^{[l]}(\mathbf{x}) := \sigma(\mathbf{W}^{[l]} f_{\boldsymbol{\theta}}^{[l-1]}(\mathbf{x}) + \mathbf{b}^{[l]})$$

with $f^{[l]} \in \mathbb{R}^{m_l}$, where $\sigma(\cdot) : \mathbb{R} \rightarrow \mathbb{R}$ is the activation function applied coordinate-wisely, with slight abuse of notation. Finally, we denote

$$f_{\boldsymbol{\theta}}(\mathbf{x}) := f(\mathbf{x}, \boldsymbol{\theta}) := f_{\boldsymbol{\theta}}^{[L]}(\mathbf{x}) := \mathbf{W}^{[L]} f_{\boldsymbol{\theta}}^{[L-1]}(\mathbf{x}) + \mathbf{b}^{[L]}, \quad (3.3)$$

and for simplicity, sometimes we may drop the subscript $\boldsymbol{\theta}$ in $f_{\boldsymbol{\theta}}^{[l]}$ for $l \in [0 : L]$.

3.2 Loss function

The set of training data is denoted by $S := \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^n$, where $\mathbf{x}_i \in \mathbb{R}^d$, $\mathbf{y}_i \in \mathbb{R}^{d'}$. Here we assume that there exists a function $f^*(\cdot) : \mathbb{R}^d \rightarrow \mathbb{R}^{d'}$ satisfying $f^*(\mathbf{x}_i) = \mathbf{y}_i$ for $i \in [n]$. We remark that this assumption helps simplify the notation in our work, however, it is not essential for our results. The empirical risk reads as

$$R_S(\boldsymbol{\theta}) := \frac{1}{n} \sum_{i=1}^n \ell(f(\mathbf{x}_i, \boldsymbol{\theta}), \mathbf{y}_i) := \mathbb{E}_S \ell(f(\mathbf{x}, \boldsymbol{\theta}), f^*(\mathbf{x})), \quad (3.4)$$

where the expectation \mathbb{E}_S is defined for any function $h(\cdot) : \mathbb{R}^d \rightarrow \mathbb{R}$ as

$$\mathbb{E}_S h(\mathbf{x}) := \frac{1}{n} \sum_{i=1}^n h(\mathbf{x}_i),$$

and the loss function $\ell(\cdot, \cdot) : \mathbb{R}^{d'} \times \mathbb{R}^{d'} \rightarrow \mathbb{R}$ in (3.4) is differentiable in both variables, and the derivative of $\ell(\cdot, \cdot)$ with respect to its first argument is denoted by $\nabla \ell(\mathbf{y}, \mathbf{y}^*)$. In this paper, we always take derivatives/gradients of $\ell(\cdot, \cdot)$ in its first argument with respect to any possible parameter.

For each $l \in [L]$, we define the error vectors as

$$\mathbf{z}_\theta^{[l]} := \nabla_{f^{[l]}} \ell$$

with $\mathbf{z}_\theta^{[l]} \in \mathbb{R}^{m_l}$, and the feature gradients as

$$\mathbf{g}_\theta^{[L]} := \mathbf{1}, \text{ and } \mathbf{g}_\theta^{[l]} := \sigma^{(1)} \left(\mathbf{W}^{[l]} \mathbf{f}_\theta^{[l-1]} + \mathbf{b}^{[l]} \right) \text{ for } l \in [L-1]$$

with $\mathbf{g}_\theta^{[l]} \in \mathbb{R}^{m_l}$, where we use $\sigma^{(1)}(\cdot)$ for the first derivative of $\sigma(\cdot)$. Moreover, we call $\mathbf{f}_\theta^{[l]}$ the feature vectors, and we denote the collections of feature vectors, feature gradients, and error vectors $\{\mathbf{z}_\theta^{[l]}\}_{l=1}^L$ respectively by

$$\mathbf{F}_\theta := \{\mathbf{f}_\theta^{[l]}\}_{l=1}^L, \mathbf{G}_\theta := \{\mathbf{g}_\theta^{[l]}\}_{l=1}^L, \mathbf{Z}_\theta := \{\mathbf{z}_\theta^{[l]}\}_{l=1}^L.$$

Moreover, using backpropagation, we can derive the following relations concerning the above quantities

$$\mathbf{z}_\theta^{[L]} = \nabla \ell, \tag{3.5a}$$

$$\mathbf{z}_\theta^{[l]} = (\mathbf{W}^{[l+1]})^\top \left(\mathbf{z}_\theta^{[l+1]} \circ \mathbf{g}_\theta^{[l+1]} \right), \quad l \in [L-1], \tag{3.5b}$$

$$\nabla_{\mathbf{W}^{[l]}} \ell = \left(\mathbf{z}_\theta^{[l]} \circ \mathbf{g}_\theta^{[l]} \right) (\mathbf{f}_\theta^{[l-1]})^\top, \quad l \in [L], \tag{3.5c}$$

$$\nabla_{\mathbf{b}^{[l]}} \ell = \mathbf{z}_\theta^{[l]} \circ \mathbf{g}_\theta^{[l]}, \quad l \in [L], \tag{3.5d}$$

where we use \circ for the Hadamard product [47] of two matrices of the same dimension.

Specifically, for simplicity of gradient computation, we define another group of error vectors for $l \in [L]$

$$\mathbf{e}_\theta^{[l]} := \mathbf{z}_\theta^{[l]} \circ \mathbf{g}_\theta^{[l]}$$

with $\mathbf{e}_\theta^{[l]} \in \mathbb{R}^{m_l}$, and we denote $\{\mathbf{e}_\theta^{[l]}\}_{l=1}^L$ by $\mathbf{E}_\theta := \{\mathbf{e}_\theta^{[l]}\}_{l=1}^L$.

Directly from relation (3.5), we obtain that

$$\mathbf{e}_\theta^{[l]} = \mathbf{z}_\theta^{[l]} \circ \mathbf{g}_\theta^{[l]} = \left((\mathbf{W}^{[l+1]})^\top \left(\mathbf{z}_\theta^{[l+1]} \circ \mathbf{g}_\theta^{[l+1]} \right) \right) \circ \mathbf{g}_\theta^{[l]} = \left((\mathbf{W}^{[l+1]})^\top \mathbf{e}_\theta^{[l+1]} \right) \circ \mathbf{g}_\theta^{[l]}. \tag{3.6}$$

3.3 Hessian

Given a scalar loss function $\ell(\cdot, \cdot) : \mathbb{R}^{d'} \times \mathbb{R}^{d'} \rightarrow \mathbb{R}$, twice differentiable in both variables, and an activation function $\sigma(\cdot) : \mathbb{R} \rightarrow \mathbb{R}$, also twice differentiable, we denote that

$$R_S(\boldsymbol{\theta}) = \mathbb{E}_S \ell(\mathbf{f}(\mathbf{x}, \boldsymbol{\theta}), \mathbf{f}^*(\mathbf{x})), \quad (3.7)$$

$$\mathbf{v}_S(\boldsymbol{\theta}) := \nabla_{\boldsymbol{\theta}} R_S(\boldsymbol{\theta}) = \mathbb{E}_S \nabla \ell(\mathbf{f}(\mathbf{x}, \boldsymbol{\theta}), \mathbf{f}^*(\mathbf{x}))^\top \nabla_{\boldsymbol{\theta}} \mathbf{f}_{\boldsymbol{\theta}}(\mathbf{x}) = \sum_{i=1}^{m_L} \mathbb{E}_S \partial_i \ell(\mathbf{f}_{\boldsymbol{\theta}}, \mathbf{f}^*) \nabla_{\boldsymbol{\theta}} (\mathbf{f}_{\boldsymbol{\theta}})_i, \quad (3.8)$$

where $\partial_i \ell(\mathbf{f}_{\boldsymbol{\theta}}, \mathbf{f}^*)$ is the i -th element of $\nabla \ell(\mathbf{f}(\mathbf{x}, \boldsymbol{\theta}), \mathbf{f}^*(\mathbf{x}))$, and $(\mathbf{f}_{\boldsymbol{\theta}})_i$ is the i -th element of vector $\mathbf{f}_{\boldsymbol{\theta}}$. Then for the Hessian matrix $\mathbf{H}_S(\boldsymbol{\theta})$, we have

$$\begin{aligned} \mathbf{H}_S(\boldsymbol{\theta}) &:= \nabla_{\boldsymbol{\theta}} \nabla_{\boldsymbol{\theta}} R_S(\boldsymbol{\theta}) \\ &= \sum_{i=1}^{m_L} \mathbb{E}_S \nabla_{\boldsymbol{\theta}} (\partial_i \ell(\mathbf{f}_{\boldsymbol{\theta}}, \mathbf{f}^*)) \nabla_{\boldsymbol{\theta}} (\mathbf{f}_{\boldsymbol{\theta}})_i + \sum_{i=1}^{m_L} \mathbb{E}_S \partial_i \ell(\mathbf{f}_{\boldsymbol{\theta}}, \mathbf{f}^*) \nabla_{\boldsymbol{\theta}} \nabla_{\boldsymbol{\theta}} ((\mathbf{f}_{\boldsymbol{\theta}})_i) \\ &= \sum_{i,j=1}^{m_L} \mathbb{E}_S \partial_{ij} \ell(\mathbf{f}_{\boldsymbol{\theta}}, \mathbf{f}^*) \nabla_{\boldsymbol{\theta}} (\mathbf{f}_{\boldsymbol{\theta}})_i (\nabla_{\boldsymbol{\theta}} (\mathbf{f}_{\boldsymbol{\theta}})_j)^\top + \sum_{i=1}^{m_L} \mathbb{E}_S \partial_i \ell(\mathbf{f}_{\boldsymbol{\theta}}, \mathbf{f}^*) \nabla_{\boldsymbol{\theta}} \nabla_{\boldsymbol{\theta}} ((\mathbf{f}_{\boldsymbol{\theta}})_i), \end{aligned}$$

where $\partial_{ij} \ell(\mathbf{f}_{\boldsymbol{\theta}}, \mathbf{f}^*)$ is the (i, j) -th element of $\nabla \nabla \ell(\mathbf{f}(\mathbf{x}, \boldsymbol{\theta}), \mathbf{f}^*(\mathbf{x}))$, with $\mathbf{H}_S(\boldsymbol{\theta}) \in \mathbb{R}^M \times \mathbb{R}^M$.

We define matrices $\mathbf{H}_S^{(1)}(\boldsymbol{\theta})$ and $\mathbf{H}_S^{(2)}(\boldsymbol{\theta})$ as follows:

$$\mathbf{H}_S^{(1)}(\boldsymbol{\theta}) := \sum_{i,j=1}^{m_L} \mathbb{E}_S \partial_{ij} \ell(\mathbf{f}_{\boldsymbol{\theta}}, \mathbf{f}^*) \nabla_{\boldsymbol{\theta}} (\mathbf{f}_{\boldsymbol{\theta}})_i (\nabla_{\boldsymbol{\theta}} (\mathbf{f}_{\boldsymbol{\theta}})_j)^\top, \quad (3.9)$$

$$\mathbf{H}_S^{(2)}(\boldsymbol{\theta}) := \sum_{i=1}^{m_L} \mathbb{E}_S \partial_i \ell(\mathbf{f}_{\boldsymbol{\theta}}, \mathbf{f}^*) \nabla_{\boldsymbol{\theta}} \nabla_{\boldsymbol{\theta}} ((\mathbf{f}_{\boldsymbol{\theta}})_i), \quad (3.10)$$

then obviously $\mathbf{H}_S^{(1)}(\boldsymbol{\theta}), \mathbf{H}_S^{(2)}(\boldsymbol{\theta}) \in \mathbb{R}^{M \times M}$, and

$$\mathbf{H}_S(\boldsymbol{\theta}) = \mathbf{H}_S^{(1)}(\boldsymbol{\theta}) + \mathbf{H}_S^{(2)}(\boldsymbol{\theta}).$$

3.4 Assumptions and conventions of notations

We begin this part by introducing several assumptions that will be used throughout this paper:

Assumptions.

(i) We choose the L -layer ($L \geq 2$) fully-connected deep neural networks (NNs) as our model.

(ii) Our training data is $S = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^n, n \in \mathbb{Z}^+$.

(iii) We use the empirical loss $R_S(\boldsymbol{\theta}) = \mathbb{E}_S \ell(\mathbf{f}_{\boldsymbol{\theta}}(\mathbf{x}), \mathbf{y})$.

(iv) Loss function $\ell(\cdot, \cdot)$ and activation function $\sigma(\cdot)$ are (weakly) differentiable. (Remark: twice differentiable is required for the computation of Hessian)

After stating out the assumptions, we would also like to introduce some conventions of notations that are frequently used in this paper in the following.

We write $\text{NN}(\{m_l\}_{l=0}^L)$ for a fully-connected L -layer network with width (m_0, \dots, m_L) , by which the tuple class of its parameters $\theta \in \text{Tuple}_{\{m_0, \dots, m_L\}}$ is determined whereas its activation $\sigma(\cdot)$ is not provided. When $\sigma(\cdot)$ is given, the output of $\text{NN}(\{m_l\}_{l=0}^L)$ is denoted by $f_\theta(x)$ with $\theta \in \text{Tuple}_{\{m_0, \dots, m_L\}}$.

Note that, if given two NNs, $\text{NN}(\{m_l\}_{l=0}^L)$ and $\text{NN}(\{m'_l\}_{l=0}^L)$, their corresponding parameters belong to different tuple classes except when $m'_l = m_l$ for all $l \in [0 : L]$. Therefore, in this work, $f_\theta(x)$ and $R_S(\theta)$ may correspond to output and loss landscape of different NNs distinguished by θ of different tuple classes.

Given two NNs, $\text{NN}(\{m_l\}_{l=0}^L)$ and $\text{NN}(\{m'_l\}_{l=0}^L)$ with $m'_0 = m_0$, $m'_L = m_L$, and $m'_l \geq m_l$ for any $l \in [L - 1]$, then for $K = \sum_{l=1}^{L-1} (m'_l - m_l) \in \mathbb{Z}^+$, we say that $\text{NN}(\{m'_l\}_{l=0}^L)$ is K -neuron **wider** than $\text{NN}(\{m_l\}_{l=0}^L)$, and conversely, $\text{NN}(\{m_l\}_{l=0}^L)$ is K -neuron **narrower** than $\text{NN}(\{m'_l\}_{l=0}^L)$.

As long as we have two NNs, $\text{NN}(\{m_l\}_{l=0}^L)$ and $\text{NN}(\{m'_l\}_{l=0}^L)$, given in the context of Definitions, Theorems, Propositions, Lemmas etc., we always assume that $\text{NN}(\{m'_l\}_{l=0}^L)$ is wider than $\text{NN}(\{m_l\}_{l=0}^L)$, i.e., $m'_0 = m_0$, $m'_L = m_L$, and $m'_l \geq m_l$ for any $l \in [L - 1]$. We also denote $M = \sum_{l=0}^{L-1} (m_l + 1)m_{l+1}$ and $M' = \sum_{l=0}^{L-1} (m'_l + 1)m'_{l+1}$, and consequently, $M' \geq M$. We denote the parameters of a narrower network by θ_{narr} , and the counterpart of a wider network by θ_{wide} . Then, given the data S , loss $\ell(\cdot, \cdot)$ and activation $\sigma(\cdot)$, the collection of critical points of narrower NN and wider NN can be found respectively and denoted by $\Theta_{\text{narr}}^c := \{\theta | \nabla_\theta R_S(\theta_{\text{narr}}) = \mathbf{0}\}$ and $\Theta_{\text{wide}}^c := \{\theta | \nabla_\theta R_S(\theta_{\text{wide}}) = \mathbf{0}\}$. Furthermore, $\mathcal{F}_{\text{narr}}^c := \{f_\theta | \theta \in \Theta_{\text{narr}}^c\}$ and $\mathcal{F}_{\text{wide}}^c := \{f_\theta | \theta \in \Theta_{\text{wide}}^c\}$ are denoted for the function spaces induced by critical points accordingly.

Finally, a neuron, say the i -th neuron in layer l , is termed a **null neuron** if its output is a constant independent of input x for any activation, i.e., $(f^{[l]})_i(\cdot) \equiv \text{Const}$ for any $\sigma(\cdot)$. Otherwise, we call this neuron an **effective neuron**.

4 Embedding Principle

In this section, we prove the Embedding Principle by constructing critical embeddings. First, we define a critical embedding operation. Then, by constructing one-step critical embeddings and their composition, we prove the Embedding Principle that critical points of the loss landscape of a narrow network can be embedded to critical affine subspaces of the loss landscape of any wider network while preserving the output function. Finally, we emphasize the importance of Embedding Principle in understanding the implicit regularization and generalization of NNs.

We begin with the concepts of embedding, affine embedding and critical embedding.

Definition 4.1 (Embedding and affine embedding). *Given an $\text{NN}(\{m_l\}_{l=0}^L)$ and $\text{NN}(\{m'_l\}_{l=0}^L)$, an **embedding** is an injective operator $\mathcal{T} : \text{Tuple}_{\{m_0, \dots, m_L\}} \rightarrow \text{Tuple}_{\{m'_0, \dots, m'_L\}}$, i.e., $\mathcal{T}(\theta_1) \neq \mathcal{T}(\theta_2)$ for $\theta_1, \theta_2 \in \text{Tuple}_{\{m_0, \dots, m_L\}}$ and $\theta_1 \neq \theta_2$. In addition, \mathcal{T} is an **affine***

embedding if $\tilde{\mathcal{T}}(\boldsymbol{\theta}) := \mathcal{T}(\boldsymbol{\theta}) - \mathcal{T}(\mathbf{0})$ is a linear operator, i.e., $\tilde{\mathcal{T}}(\boldsymbol{\theta}_1) + \tilde{\mathcal{T}}(\boldsymbol{\theta}_2) = \tilde{\mathcal{T}}(\boldsymbol{\theta}_1 + \boldsymbol{\theta}_2)$ and $\tilde{\mathcal{T}}(\beta\boldsymbol{\theta}) = \beta\tilde{\mathcal{T}}(\boldsymbol{\theta})$ for any $\boldsymbol{\theta}, \boldsymbol{\theta}_1, \boldsymbol{\theta}_2 \in \text{Tuple}_{\{m_0, \dots, m_L\}}$ and $\beta \in \mathbb{R}$.

Remark 4.1. For any given affine embedding \mathcal{T} , it is associated with a matrix $\mathbf{A} \in \mathbb{R}^{M' \times M}$ and a vector $\mathbf{c} \in \mathbb{R}^{M'}$ such that $\text{vec}(\mathcal{T}(\boldsymbol{\theta})) = \text{Avec}(\boldsymbol{\theta}) + \mathbf{c}$, where $M = \sum_{l=0}^{L-1} (m_l + 1)m_{l+1}$ and $M' = \sum_{l=0}^{L-1} (m'_l + 1)m'_{l+1}$. As noted before, we do not distinguish tuple $\boldsymbol{\theta}$ from its vectorization $\text{vec}(\boldsymbol{\theta})$ in the following. Hence, $\mathcal{T}(\boldsymbol{\theta}) = \mathbf{A}\boldsymbol{\theta} + \mathbf{c}$.

Definition 4.2 (Critical embedding). Given an $\text{NN}(\{m_l\}_{l=0}^L)$ and $\text{NN}(\{m'_l\}_{l=0}^L)$, a **critical embedding** is an affine embedding $\mathcal{T} : \text{Tuple}_{\{m_0, \dots, m_L\}} \rightarrow \text{Tuple}_{\{m'_0, \dots, m'_L\}}$, which maps any set of its network parameters $\boldsymbol{\theta}_{\text{narr}} \in \text{Tuple}_{\{m_0, \dots, m_L\}}$ to that of a wider NN $\boldsymbol{\theta}_{\text{wide}} = \mathcal{T}(\boldsymbol{\theta}_{\text{narr}}) \in \text{Tuple}_{\{m'_0, \dots, m'_L\}}$ satisfying that: For any given data S , loss function $\ell(\cdot, \cdot)$, activation function $\sigma(\cdot)$,

(i) **output preserving:** $f_{\boldsymbol{\theta}_{\text{narr}}}(\mathbf{x}) = f_{\boldsymbol{\theta}_{\text{wide}}}(\mathbf{x})$ for any $\mathbf{x} \in \mathbb{R}^d$;

(ii) **representation preserving:**

$$\begin{aligned} & \text{span} \left\{ \left\{ \left(f_{\boldsymbol{\theta}_{\text{narr}}}^{[l]}(\cdot) \right)_j \right\}_{j \in [m_l]} \cup \{1\} \right\} \\ &= \text{span} \left\{ \left\{ \left(f_{\boldsymbol{\theta}_{\text{wide}}}^{[l]}(\cdot) \right)_{j'} \right\}_{j' \in [m'_l]} \cup \{1\} \right\}, \quad \text{for any } l \in [L], \end{aligned}$$

where $\{f_{\boldsymbol{\theta}_{\text{narr}}}^{[l]}\}_{l=1}^L$ and $\{f_{\boldsymbol{\theta}_{\text{wide}}}^{[l]}\}_{l=1}^L$ are feature vectors of $\text{NN}(\{m_l\}_{l=0}^L)$ and $\text{NN}(\{m'_l\}_{l=0}^L)$, and $1 : \mathbb{R}^d \rightarrow \mathbb{R}$ is the constant function, i.e., $1(\cdot) \equiv 1$;

(iii) **criticality preserving:** If $\boldsymbol{\theta}_{\text{narr}}$ is a critical point of $R_S(\boldsymbol{\theta})$, i.e., $\nabla_{\boldsymbol{\theta}} R_S(\boldsymbol{\theta}_{\text{narr}}) = \mathbf{0}$, then $\boldsymbol{\theta}_{\text{wide}}$ is also a critical point of $R_S(\boldsymbol{\theta})$, i.e., $\nabla_{\boldsymbol{\theta}} R_S(\boldsymbol{\theta}_{\text{wide}}) = \mathbf{0}$.

Specifically, if an embedding is a critical embedding with $\text{NN}(\{m'_l\}_{l=0}^L)$ one-neuron wider than $\text{NN}(\{m_l\}_{l=0}^L)$, we call it **one-step critical embedding**.

4.1 One-step critical embedding

In this subsection, we introduce two types of one-step critical embeddings. we start with the definition of one-step null embedding. Intuitively, a one-step null embedding adds a null neuron to the NN, whose output is a constant independent of input \mathbf{x} for any activation.

Definition 4.3 (One-step null embedding). Given an $\text{NN}(\{m_l\}_{l=0}^L)$ and its parameter $\boldsymbol{\theta} = (\mathbf{W}^{[1]}, \mathbf{b}^{[1]}, \dots, \mathbf{W}^{[L]}, \mathbf{b}^{[L]}) \in \text{Tuple}_{\{m_0, \dots, m_L\}}$, then for any $l \in [L-1]$, we define the operators $\mathcal{T}_{l,0}$ and $\mathcal{V}_{l,0}$ applying on $\boldsymbol{\theta}$ as follows

$$\begin{aligned} \mathcal{T}_{l,0}(\boldsymbol{\theta})|_k &= \boldsymbol{\theta}|_k, \quad k \neq l, l+1, \\ \mathcal{T}_{l,0}(\boldsymbol{\theta})|_l &= \left(\left[\begin{array}{c} \mathbf{W}^{[l]} \\ \mathbf{0}_{1 \times m_{l-1}} \end{array} \right], \left[\begin{array}{c} \mathbf{b}^{[l]} \\ 0 \end{array} \right] \right), \end{aligned}$$

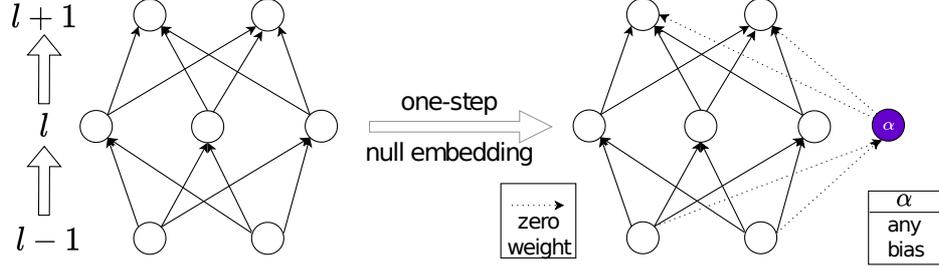


Figure 4.1: Illustration of one-step null embedding. The purple neuron is added with both input and output weights as zero and an arbitrary bias α .

$$\begin{aligned}\mathcal{T}_{l,0}(\boldsymbol{\theta})|_{l+1} &= \left(\left[\mathbf{W}^{[l+1]}, \mathbf{0}_{m_{l+1} \times 1} \right], \mathbf{b}^{[l+1]} \right), \\ \mathcal{V}_{l,0}(\boldsymbol{\theta})|_k &= (\mathbf{0}_{m_k \times m_{k-1}}, \mathbf{0}_{m_k \times 1}), \quad k \neq l, l+1, \\ \mathcal{V}_{l,0}(\boldsymbol{\theta})|_l &= \left(\mathbf{0}_{(m_l+1) \times m_{l-1}}, \begin{bmatrix} \mathbf{0}_{m_l \times 1} \\ 1 \end{bmatrix} \right), \\ \mathcal{V}_{l,0}(\boldsymbol{\theta})|_{l+1} &= (\mathbf{0}_{m_{l+1} \times (m_l+1)}, \mathbf{0}_{m_{l+1} \times 1}).\end{aligned}$$

We define **one-step null embedding** $\mathcal{T}_{l,0}^\alpha$ as: For any $\boldsymbol{\theta} \in \text{Tuple}_{\{m_0, \dots, m_L\}}$,

$$\mathcal{T}_{l,0}^\alpha(\boldsymbol{\theta}) = (\mathcal{T}_{l,0} + \alpha \mathcal{V}_{l,0})(\boldsymbol{\theta}).$$

Note that the neuron added by the above one-step null embedding has zero output weights, zero input weights and an arbitrary bias α . An illustration is shown in Fig. 4.1.

We proceed to the definition of one-step splitting embedding.

Definition 4.4 (One-step splitting embedding). Given an $\text{NN}(\{m_l\}_{l=0}^L)$ and its parameter $\boldsymbol{\theta} = (\mathbf{W}^{[1]}, \mathbf{b}^{[1]}, \dots, \mathbf{W}^{[L]}, \mathbf{b}^{[L]}) \in \text{Tuple}_{\{m_0, \dots, m_L\}}$, then for any $l \in [L-1]$ and $s \in [m_l]$, we define the operators $\mathcal{T}_{l,s}$ and $\mathcal{V}_{l,s}$ applying on $\boldsymbol{\theta}$ as follows

$$\begin{aligned}\mathcal{T}_{l,s}(\boldsymbol{\theta})|_k &= \boldsymbol{\theta}|_k, \quad k \neq l, l+1, \\ \mathcal{T}_{l,s}(\boldsymbol{\theta})|_l &= \left(\begin{bmatrix} \mathbf{W}^{[l]} \\ \mathbf{W}_{s, [1:m_{l-1}]}^{[l]} \end{bmatrix}, \begin{bmatrix} \mathbf{b}^{[l]} \\ \mathbf{b}_s^{[l]} \end{bmatrix} \right), \\ \mathcal{T}_{l,s}(\boldsymbol{\theta})|_{l+1} &= \left(\left[\mathbf{W}^{[l+1]}, \mathbf{0}_{m_{l+1} \times 1} \right], \mathbf{b}^{[l+1]} \right), \\ \mathcal{V}_{l,s}(\boldsymbol{\theta})|_k &= (\mathbf{0}_{m_k \times m_{k-1}}, \mathbf{0}_{m_k \times 1}), \quad k \neq l, l+1, \\ \mathcal{V}_{l,s}(\boldsymbol{\theta})|_l &= (\mathbf{0}_{(m_l+1) \times m_{l-1}}, \mathbf{0}_{(m_l+1) \times 1}), \\ \mathcal{V}_{l,s}(\boldsymbol{\theta})|_{l+1} &= \left(\left[\mathbf{0}_{m_{l+1} \times (s-1)}, -\mathbf{W}_{[1:m_{l+1}], s}^{[l+1]}, \mathbf{0}_{m_{l+1} \times (m_l-s)}, \mathbf{W}_{[1:m_{l+1}], s}^{[l+1]} \right], \mathbf{0}_{m_{l+1} \times 1} \right).\end{aligned}$$

We define **one-step splitting embedding** $\mathcal{T}_{l,s}^\alpha$ as: For any $\boldsymbol{\theta} \in \text{Tuple}_{\{m_0, \dots, m_L\}}$,

$$\mathcal{T}_{l,s}^\alpha(\boldsymbol{\theta}) = (\mathcal{T}_{l,s} + \alpha \mathcal{V}_{l,s})(\boldsymbol{\theta}).$$

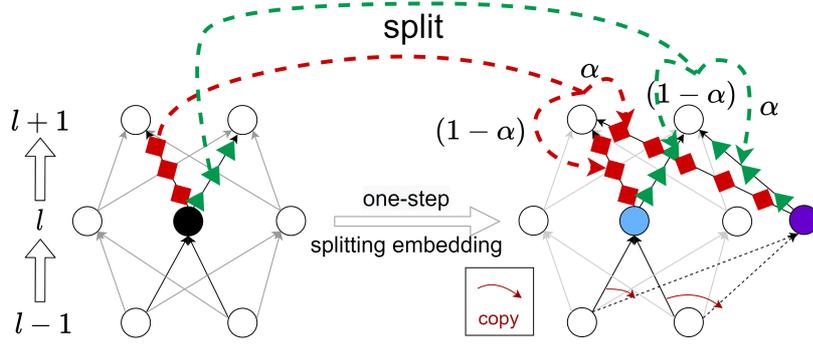


Figure 4.2: Illustration of one-step splitting embedding. The black neuron in the left network is split into the blue and purple neurons in the right network. The red (green) output weight of the black neuron in the left net is split into two red (green) weights in the right net with ratio $(1 - \alpha)$ and α , respectively. This is also illustrated in [24].

An illustration of one-step splitting method is shown in Fig. 4.2.

Remark 4.2. The parameters $\mathcal{T}_{l,0}^\alpha(\theta)$ and $\mathcal{T}_{l,s}^\alpha(\theta)$ correspond to a L -layer NN with width $(m_0, \dots, m_{l-1}, m_l + 1, m_{l+1}, \dots, m_L)$ since $\mathcal{T}_{l,0}^\alpha$ and $\mathcal{T}_{l,s}^\alpha$ are one-step embeddings.

Remark 4.3. We observe that $\mathcal{T}_{l,0}^\alpha$ and $\mathcal{T}_{l,s}^\alpha$ can be applied on the neural network parameter θ of any given NN($\{m_l\}_{l=0}^L$) of proper depth and width, hence the domain of $\mathcal{T}_{l,0}^\alpha$ and $\mathcal{T}_{l,s}^\alpha$ is not limited to a specific Tuple $_{\{m_0, \dots, m_L\}}$. Instead, the extended domain of $\mathcal{T}_{l,0}^\alpha$ and $\mathcal{T}_{l,s}^\alpha$ assembles all possible tuple class, i.e., if we denote the extended domain of $\mathcal{T}_{l,0}^\alpha$ by $\mathcal{D}_{l,0}$, then

$$\mathcal{D}_{l,0} := \bigsqcup_{l < L} \text{Tuple}_{\{n_0, \dots, n_L\}}, \quad (4.1)$$

and if we denote the extended domain of $\mathcal{T}_{l,s}^\alpha$ by $\mathcal{D}_{l,s}$, then

$$\mathcal{D}_{l,s} := \bigsqcup_{l < L, s \leq n_l} \text{Tuple}_{\{n_0, \dots, n_L\}}, \quad (4.2)$$

where \bigsqcup refers to the disjoint union of different tuple classes.

By the above remark, we may extend $\mathcal{T}_{l,0}^\alpha$ and $\mathcal{T}_{l,s}^\alpha$ to their extended domains, and we identify

$$\mathcal{T}_{l,0}^\alpha : \mathcal{D}_{l,0} \rightarrow \mathcal{D}_{l,0}, \quad \mathcal{T}_{l,s}^\alpha : \mathcal{D}_{l,s} \rightarrow \mathcal{D}_{l,s},$$

with their restrictions on $\text{Tuple}_{\{m_0, \dots, m_L\}}$ for some given NN($\{m_l\}_{l=0}^L$).

Theorem 4.1. *One-step null embedding and one-step splitting embedding are critical embeddings.*

In order to prove Theorem 4.1, we need Lemma 4.1 and Lemma 4.2, where Lemma 4.2 has already been presented previously in our conference paper [24, Lemma 1].

Lemma 4.1. For any one-step null embedding $\mathcal{T}_{l,0}^\alpha$, given any $\text{NN}(\{m_l\}_{l=0}^L)$ and its parameters $\boldsymbol{\theta}_{\text{narr}} \in \text{Tuple}_{\{m_0, \dots, m_L\}}$ with $\text{Tuple}_{\{m_0, \dots, m_L\}} \in \mathcal{D}_{l,0}$, we have $\boldsymbol{\theta}_{\text{wide}} := \mathcal{T}_{l,0}^\alpha(\boldsymbol{\theta}_{\text{narr}})$ satisfies the following conditions: given any data S , loss $\ell(\cdot, \cdot)$ and activation $\sigma(\cdot)$, for any $l \in [L-1]$,

(i) feature vectors in

$$\mathbf{F}_{\boldsymbol{\theta}_{\text{wide}}}: \mathbf{f}_{\boldsymbol{\theta}_{\text{wide}}}^{[l']} = \mathbf{f}_{\boldsymbol{\theta}_{\text{narr}}}^{[l']}, \text{ for } l' \in [L] \text{ and } l' \neq l, \mathbf{f}_{\boldsymbol{\theta}_{\text{wide}}}^{[l]} = \left[(\mathbf{f}_{\boldsymbol{\theta}_{\text{narr}}}^{[l]})^\top, \sigma(\alpha) \right]^\top;$$

(ii) feature gradients in

$$\mathbf{G}_{\boldsymbol{\theta}_{\text{wide}}}: \mathbf{g}_{\boldsymbol{\theta}_{\text{wide}}}^{[l']} = \mathbf{g}_{\boldsymbol{\theta}_{\text{narr}}}^{[l']}, \text{ for } l' \in [L] \text{ and } l' \neq l, \mathbf{g}_{\boldsymbol{\theta}_{\text{wide}}}^{[l]} = \left[(\mathbf{g}_{\boldsymbol{\theta}_{\text{narr}}}^{[l]})^\top, \sigma^{(1)}(\alpha) \right]^\top;$$

(iii) error vectors in

$$\mathbf{Z}_{\boldsymbol{\theta}_{\text{wide}}}: \mathbf{z}_{\boldsymbol{\theta}_{\text{wide}}}^{[l']} = \mathbf{z}_{\boldsymbol{\theta}_{\text{narr}}}^{[l']}, \text{ for } l' \in [L] \text{ and } l' \neq l, \mathbf{z}_{\boldsymbol{\theta}_{\text{wide}}}^{[l]} = \left[\mathbf{z}_{\boldsymbol{\theta}_{\text{narr}}}^{[l]}, 0 \right]^\top;$$

(iv) $\mathcal{T}_{l,0}^\alpha$ is injective for all α ;

(v) $\mathcal{T}_{l,0}^\alpha$ is an affine embedding for all α .

Lemma 4.2 (Lemma 1 in [24]). For any one-step splitting embedding $\mathcal{T}_{l,s}^\alpha$, given any $\text{NN}(\{m_l\}_{l=0}^L)$ and its parameters $\boldsymbol{\theta}_{\text{narr}} \in \text{Tuple}_{\{m_0, \dots, m_L\}}$ with $\text{Tuple}_{\{m_0, \dots, m_L\}} \in \mathcal{D}_{l,s}$, we have $\boldsymbol{\theta}_{\text{wide}} := \mathcal{T}_{l,s}^\alpha(\boldsymbol{\theta}_{\text{narr}})$ satisfies the following conditions: given any data S , loss $\ell(\cdot, \cdot)$ and activation $\sigma(\cdot)$, for any $l \in [L-1]$,

(i) feature vectors in

$$\mathbf{F}_{\boldsymbol{\theta}_{\text{wide}}}: \mathbf{f}_{\boldsymbol{\theta}_{\text{wide}}}^{[l']} = \mathbf{f}_{\boldsymbol{\theta}_{\text{narr}}}^{[l']}, \text{ for } l' \in [L] \text{ and } l' \neq l, \mathbf{f}_{\boldsymbol{\theta}_{\text{wide}}}^{[l]} = \left[(\mathbf{f}_{\boldsymbol{\theta}_{\text{narr}}}^{[l]})^\top, (\mathbf{f}_{\boldsymbol{\theta}_{\text{narr}}}^{[l]})_s \right]^\top;$$

(ii) feature gradients in

$$\mathbf{G}_{\boldsymbol{\theta}_{\text{wide}}}: \mathbf{g}_{\boldsymbol{\theta}_{\text{wide}}}^{[l']} = \mathbf{g}_{\boldsymbol{\theta}_{\text{narr}}}^{[l']}, \text{ for } l' \in [L] \text{ and } l' \neq l, \mathbf{g}_{\boldsymbol{\theta}_{\text{wide}}}^{[l]} = \left[(\mathbf{g}_{\boldsymbol{\theta}_{\text{narr}}}^{[l]})^\top, (\mathbf{g}_{\boldsymbol{\theta}_{\text{narr}}}^{[l]})_s \right]^\top;$$

(iii) error vectors in

$$\begin{aligned} \mathbf{Z}_{\boldsymbol{\theta}_{\text{wide}}}: \mathbf{z}_{\boldsymbol{\theta}_{\text{wide}}}^{[l']} &= \mathbf{z}_{\boldsymbol{\theta}_{\text{narr}}}^{[l']}, \text{ for } l' \in [L] \text{ and } l' \neq l, \\ \mathbf{z}_{\boldsymbol{\theta}_{\text{wide}}}^{[l]} &= \left[\left(\mathbf{z}_{\boldsymbol{\theta}_{\text{narr}}}^{[l]} \right)_{[1:s-1]}^\top, (1-\alpha)(\mathbf{z}_{\boldsymbol{\theta}_{\text{narr}}}^{[l]})_s, \left(\mathbf{z}_{\boldsymbol{\theta}_{\text{narr}}}^{[l]} \right)_{[s+1:m_l]}^\top, \alpha(\mathbf{z}_{\boldsymbol{\theta}_{\text{narr}}}^{[l]})_s \right]^\top; \end{aligned}$$

(iv) $\mathcal{T}_{l,s}^\alpha$ is injective for all α .

(v) $\mathcal{T}_{l,s}^\alpha$ is an affine embedding for all α .

Directly from Lemma 4.1 and Lemma 4.2, we obtain that both one-step null embedding (Proposition 4.1) and one-step splitting embedding (Proposition 4.2) satisfy the property of output preserving and representation preserving, and all we need is to check the property of criticality preserving. We remark that Proposition 4.2 has also been presented previously in our conference paper [24, Proposition 1].

Proposition 4.1. *For any one-step null embedding $\mathcal{T}_{l,0}^\alpha$, given any $\text{NN}(\{m_l\}_{l=0}^L)$ and its parameters $\theta_{\text{narr}} \in \text{Tuple}_{\{m_0, \dots, m_L\}}$ with $\text{Tuple}_{\{m_0, \dots, m_L\}} \in \mathcal{D}_{l,0}$, we have $\theta_{\text{wide}} := \mathcal{T}_{l,0}^\alpha(\theta_{\text{narr}})$ satisfies the following conditions: given any data S , loss $\ell(\cdot, \cdot)$ and activation $\sigma(\cdot)$, if $\nabla_{\theta} R_S(\theta_{\text{narr}}) = \mathbf{0}$, then $\nabla_{\theta} R_S(\theta_{\text{wide}}) = \mathbf{0}$.*

Proposition 4.2. *For any one-step splitting embedding $\mathcal{T}_{l,s}^\alpha$, given any $\text{NN}(\{m_l\}_{l=0}^L)$ and its parameters $\theta_{\text{narr}} \in \text{Tuple}_{\{m_0, \dots, m_L\}}$ with $\text{Tuple}_{\{m_0, \dots, m_L\}} \in \mathcal{D}_{l,s}$, we have $\theta_{\text{wide}} := \mathcal{T}_{l,s}^\alpha(\theta_{\text{narr}})$ satisfies the following conditions: given any data S , loss $\ell(\cdot, \cdot)$ and activation $\sigma(\cdot)$, if $\nabla_{\theta} R_S(\theta_{\text{narr}}) = \mathbf{0}$, then $\nabla_{\theta} R_S(\theta_{\text{wide}}) = \mathbf{0}$.*

Combining altogether Lemma 4.1, Lemma 4.2, Proposition 4.1 and Proposition 4.2, we finish our proof for Theorem 4.1.

4.2 Composition of one-step embeddings and the Embedding Principle

We firstly define the composition of two embeddings, which readily leads to the multi-step embedding operation. (A formal definition of the composition of two embeddings can be found in Appendix C.)

Definition 4.5 (Composition of two embeddings, Informal). *For any two embeddings \mathcal{T} and \mathcal{T}' , for any θ in the domain of \mathcal{T} , the operator $\mathcal{T}'\mathcal{T}$ defined as $\mathcal{T}'\mathcal{T}(\theta) := \mathcal{T}'(\mathcal{T}(\theta))$ is also an embedding, and we term $\mathcal{T}'\mathcal{T}$ the composition of \mathcal{T}' and \mathcal{T} .*

For simplicity, for any $K \in \mathbb{Z}^+$, $K \geq 2$, we denote hereafter $\prod_{l=1}^K \mathcal{T}_l := \mathcal{T}_K \cdots \mathcal{T}_1$ as the composition of K individual embeddings $\{\mathcal{T}_l\}_{l=1}^K$. For $K = 1$, $\prod_{l=1}^1 \mathcal{T}_l := \mathcal{T}_1$.

Definition 4.6 (K -step (Multi-step) composition embedding). *Suppose we have two vectors $\mathbf{l} = (l_k)_{k=1}^K$, $l_k \in [L-1]$, $\boldsymbol{\alpha} = (\alpha_k)_{k=1}^K \subset \mathbb{R}^K$, and a sequence $\{m_l^{(0)}\}_{l=1}^L$ with $m_l^{(0)} := m_l$ for $l \in [L]$. Then given an $\text{NN}(\{m_l\}_{l=0}^L)$ and its parameters θ , a K -step composition embedding, $\mathcal{T} : \text{Tuple}_{\{m_0, \dots, m_L\}} \rightarrow \text{Tuple}_{\{m'_0, \dots, m'_L\}}$ with $K = \sum_{l=1}^{L-1} m'_j - \sum_{l=1}^{L-1} m_l$, is defined recursively by the composition of K one-step null embeddings or one-step splitting embeddings.*

Formally speaking, a K -step composition embedding $\mathcal{T}_{l,s}^\alpha$ is defined recursively as follows:

For $n = 1$, choose $s_1 \in [m_{l_1}^{(0)}] \cup \{0\}$, then

$$\mathcal{T}_{l,s}^{\boldsymbol{\alpha},(1)} := \mathcal{T}_{l_1, s_1}^{\alpha_1} = \mathcal{T}_{l_1, s_1} + \alpha_1 \mathcal{V}_{l_1, s_1}.$$

Update the sequence from $\{m_l^{(0)}\}_{l=1}^L$ to $\{m_l^{(1)}\}_{l=1}^L$ following: $m_l^{(1)} = m_l^{(0)}$ for $l \in [L-1] \setminus \{l_1\}$, and $m_{l_1}^{(1)} = m_{l_1}^{(0)} + 1$;

Then inductively, for $n = k$, choose $s_k \in [m_{l_k}^{(k-1)}] \cup \{0\}$, then

$$\mathcal{T}_{l,s}^{\alpha,(k)} := \mathcal{T}_{l_k, s_k}^{\alpha_k} \mathcal{T}_{l,s}^{\alpha,(k-1)}.$$

Update the sequence from $\{m_l^{(k-1)}\}_{l=1}^L$ to $\{m_l^{(k)}\}_{l=1}^L$ following: $m_l^{(k)} = m_l^{(k-1)}$ for $l \in [L - 1] \setminus \{l_k\}$, and $m_{l_k}^{(k)} = m_{l_k}^{(k-1)} + 1$.

Finally, $\mathcal{T} := \mathcal{T}_{l,s}^{\alpha} := \mathcal{T}_{l,s}^{\alpha,(K)}$.

Remark 4.4. For each $i \in [K]$, $\mathcal{T}_{l_i, s_i}^{\alpha_i}$ is regarded as its restriction on the tuple class $\text{Tuple}_{\{m_0^{(i-1)}, \dots, m_L^{(i-1)}\}} \in \mathcal{D}_{l_i, s_i}$, hence

$$\mathcal{T}_{l_i, s_i}^{\alpha_i} : \text{Tuple}_{\{m_0^{(i-1)}, \dots, m_L^{(i-1)}\}} \rightarrow \text{Tuple}_{\{m_0^{(i)}, \dots, m_L^{(i)}\}}.$$

Theorem 4.2. A K -step composition embedding is a critical embedding.

The composition of one-step embeddings renders a feasible method to embedding any critical point of the loss landscape of a narrow NN to a critical point with the same output function of any wider NN, therefore, we have the following Embedding Principle.

Theorem 4.3 (Embedding Principle). *Given any NN and any K -neuron wider NN, there exists a K -step composition embedding \mathcal{T} satisfying that: For any given data S , loss function $\ell(\cdot, \cdot)$, activation function $\sigma(\cdot)$, given any critical point θ_{narr}^c of the narrower NN, $\theta_{\text{wide}}^c := \mathcal{T}(\theta_{\text{narr}}^c)$ is still a critical point of the K -neuron wider NN with the same output function, i.e., $f_{\theta_{\text{narr}}^c} = f_{\theta_{\text{wide}}^c}$.*

Proof. Existence of a K -step composition embedding \mathcal{T} can be seen from the K -step construction given in Definition 4.6, and we finish the proof. \square

Moreover, we obtain a corollary from Theorem 4.3 stating the Embedding Principle for the critical functions.

Corollary 4.1 (Embedding Principle of critical functions). *Given any NN and any wider NN, for any given data S , loss function $\ell(\cdot, \cdot)$ and activation function $\sigma(\cdot)$,*

$$\mathcal{F}_{\text{narr}}^c \subset \mathcal{F}_{\text{wide}}^c, \quad (4.3)$$

where $\mathcal{F}_{\text{narr}}^c = \{f_{\theta} | \theta \in \Theta_{\text{narr}}^c\}$ and $\mathcal{F}_{\text{wide}}^c = \{f_{\theta} | \theta \in \Theta_{\text{wide}}^c\}$ are the sets of critical functions.

4.3 Importance of Embedding Principle

Mathematically speaking, Embedding Principle is a natural result of an embedding operation that preserves output function and criticality. However, we must emphasize the importance of stating and proving it explicitly in this work. For a long time, researchers study the loss landscape of NN from an optimization perspective focusing specifically on its property in the parameter space. Lots of works make effort in tackling problems like whether bad local minima exist, whether local minima are also global minima and whether

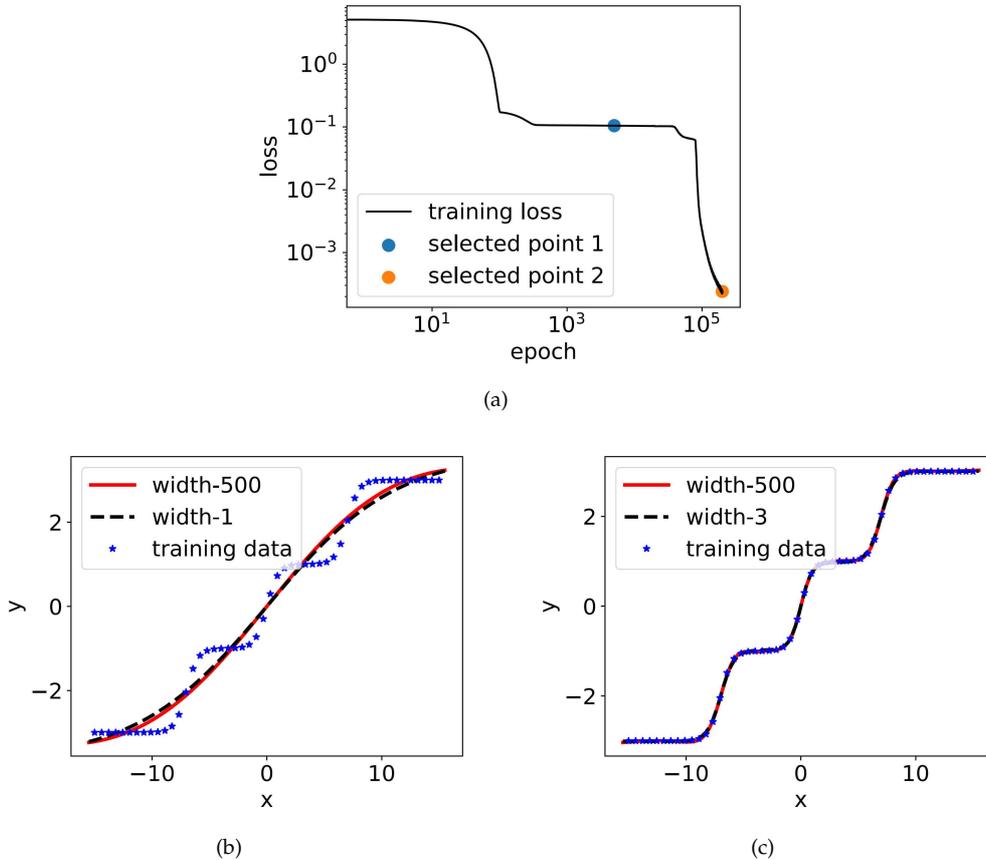


Figure 4.3: (a) The training loss of a two-layer tanh neural network with 500 hidden neurons. (b, c) Red solid: the DNN output at a training step, where the blue dot and the orange dot in (a) corresponds to (b) and (c), respectively; Black dashed: the output of the global minimum of the width-1 NN in (b) and the width-3 NN in (c), respectively; Blue dots: training data. This is also illustrated in [24].

all saddle points are strict-saddle points, etc. However, because the loss landscape of NN also has profound impact on its implicit regularization and generalization performance, it is important to look into the loss landscape from the perspective of function spaces.

Motivated by the phenomena of Frequency Principle [14–18] and condensation [11], we are very interested in the question of what are the critical functions of an NN loss landscape that may attract the training trajectory in the function space. Specifically, we care about whether there are “simple” critical functions in wide NNs that may implicitly regularize the training to help avoid overfitting. For example, in our experiments shown in Fig. 4.3, we clearly observe that the training of a width-500 two-layer tanh-NN in fitting 50 data points experiences two stages. At the first stage, it learns an output function close to the best fitting of the width-1 tanh-NN and stays for a while, seemingly that it encounters a saddle point. At the second stage, it converges to an output function close

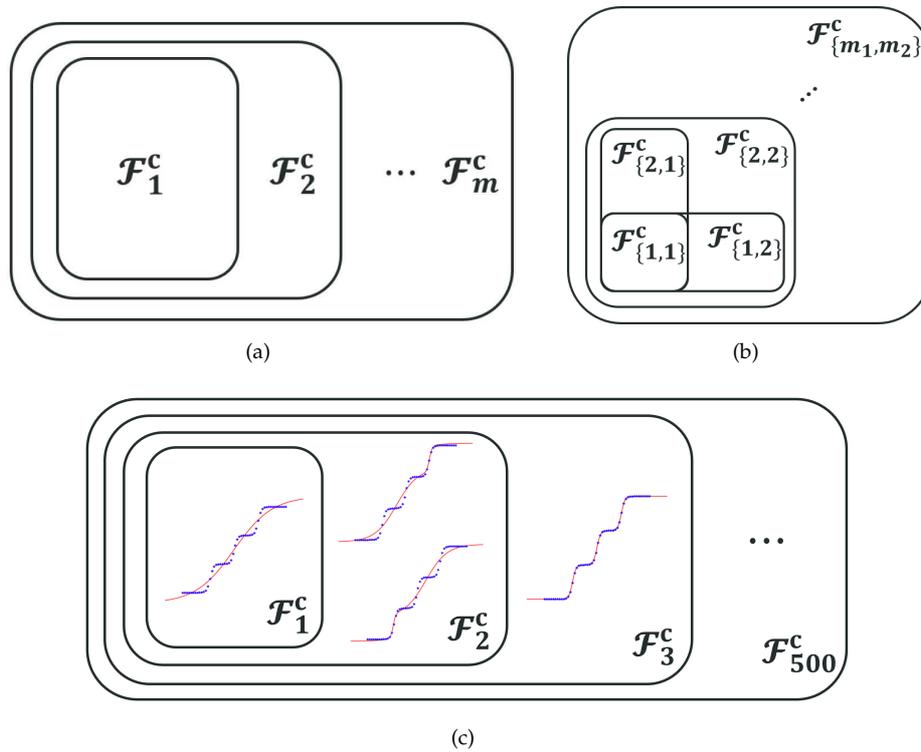


Figure 4.4: Illustration of the Embedding Principle for the critical functions of (a) two-layer width- m NN and (b) three-layer width- $\{m_1, m_2\}$ NN. (c) Some critical functions of width-500 tanh-NN predicted by the Embedding Principle for experiments in Fig. 4.3.

to the best fitting of the width-3 tanh-NN, which interpolates all the data points. Clearly, the complexity of the output function of the width-500 NN gradually increases during the training, leading to a non-overfitting interpolation of data despite of possessing overfitting capability. However, before studying the universality of such training behavior clearly relevant to generalization, it is important to have a theoretical answer to whether such “simple” critical functions always exist even in very wide NNs.

By stating and proving the Embedding Principle explicitly, we provide a clear answer of “YES” to above question. Moreover, we unravel the exact meaning of “simple” critical functions—critical functions of narrower networks. An illustration of critical functions of a two-layer width- m NN and a three-layer width- $\{m_1, m_2\}$ NN predicted by the Embedding Principle are shown in Fig. 4.4(a), (b), respectively. Critical functions of the width-500 tanh-NN for experiments in Fig. 4.3 are illustrated in Fig. 4.4(c). Note that, by Frequency Principle [14–18], we consider functions dominated by low frequency components as “simple” functions that training of an NN is implicitly biased to. Here, by Embedding Principle, “simple” functions are those critical functions of narrow NNs, which signify a hierarchical structure of fittings of training data with different complexities indicated by the narrowest NN a critical function belongs to. Combining with empirical observations of Frequency

Principle and condensation, we conjecture that nonlinear training of an NN is implicitly biased towards these “simple” critical functions. It is clearly important to look further into this conjecture in the future works.

5 General compatible critical embedding

Understanding the critical points/manifolds and their geometry of NN loss landscape is an important open problem in deep learning theory. To further advance our study along this direction, the Embedding Principle proven above inspires us to treat NN loss landscape of different widths and their critical points/manifolds as a unified object linked by the critical embeddings. These explicit critical submanifolds embedded from narrower NNs are subsets of the critical point set of a given NN. The study of them could provide lower bounds of important geometric properties of the critical point set such as dimensions of critical submanifolds (associated to the degeneracy of critical points). Clearly, the more critical embeddings we discover, the larger embedded critical submanifolds we uncover explicitly, and the tighter lower bound estimation we may obtain.

To define a general compatible embedding from a narrow NN to a wide NN, we first define mappings that establish the relation between neuron indices of the narrow NN and neuron indices of the wide NN.

Definition 5.1 (Pull-back index mapping and total index mapping). *Given an NN($\{m_l\}_{l=0}^L$) and a wider NN($\{m'_l\}_{l=0}^L$), a **pull-back index mapping** $\mathcal{I} := \{\mathcal{I}_l\}_{l=0}^L$ from the wider NN to the narrower NN is defined as follows: For any fixed $l \in [L-1]$, \mathcal{I}_l maps a neuron index $s' \in [m'_l]$ of the wider NN($\{m'_l\}_{l=0}^L$) in layer l to a neuron index $s \in [m_l] \sqcup \{0\}$ in the same layer of the narrower NN($\{m_l\}_{l=0}^L$). As for the case of $l = 0$ and $l = L$, \mathcal{I}_0 and \mathcal{I}_L are always the identity maps since their indices are fixed once data is given with $m'_0 = m_0 = d$ and $m'_L = m_L = d'$. To sum up*

$$\mathcal{I}_0 : [m'_0] \rightarrow [m_0], \quad \mathcal{I}_l : [m'_l] \rightarrow [m_l] \sqcup \{0\} \text{ for } l \in [L-1], \quad \mathcal{I}_L : [m'_L] \rightarrow [m_L]. \quad (5.1)$$

Moreover, for any nonzero index $s \neq 0$, if $\mathcal{I}_l^{-1}(s) \neq \emptyset$ for all $l \in [0:L]$, we say that $\mathcal{I} = \{\mathcal{I}_l\}_{l=0}^L$ is a **total (pull-back) index mapping**.

Lemma 5.1. *For any affine embedding $\mathcal{T} : \text{Tuple}_{\{m_0, \dots, m_L\}} \rightarrow \text{Tuple}_{\{m'_0, \dots, m'_L\}}$ satisfying the output preserving property, if there exists a total index mapping $\mathcal{I} = \{\mathcal{I}_l\}_{l=0}^L$ from NN($\{m'_l\}_{l=0}^L$) to NN($\{m_l\}_{l=0}^L$) and auxiliary variables $\beta = \{\beta_j^{[l]} \in \mathbb{R} \mid l \in [0:L], j \in [m'_l] \setminus \mathcal{I}_l^{-1}(0)\}$, such that for any given neuron belonging to NN($\{m'_l\}_{l=0}^L$), located in layer l with index j , the following two statements hold:*

$$(i) \text{ If } \mathcal{I}_l(j) \neq 0, (\mathbf{f}_{\theta_{\text{wide}}}^{[l]})_j = (\mathbf{f}_{\theta_{\text{narr}}}^{[l]})_{\mathcal{I}_l(j)} \text{ and } (\mathbf{e}_{\theta_{\text{wide}}}^{[l]})_j = \beta_j^{[l]} (\mathbf{e}_{\theta_{\text{narr}}}^{[l]})_{\mathcal{I}_l(j)},$$

$$(ii) \text{ If } \mathcal{I}_l(j) = 0, (\mathbf{f}_{\theta_{\text{wide}}}^{[l]})_j = \text{Const} \text{ and } (\mathbf{e}_{\theta_{\text{wide}}}^{[l]})_j = 0,$$

then \mathcal{T} is a critical embedding.

Next, we propose a general compatible embedding method, where above embedding methods including one-step embedding and K-step composition embedding are its special cases.

Definition 5.2 (General compatible embedding). *Given an $\text{NN}(\{m_l\}_{l=0}^L)$ and a wider $\text{NN}(\{m'_l\}_{l=0}^L)$, then for any total index mapping $\mathcal{I} = \{\mathcal{I}_l\}_{l=0}^L$ from $\text{NN}(\{m'_l\}_{l=0}^L)$ to $\text{NN}(\{m_l\}_{l=0}^L)$, and for any tuple $\alpha := \{\alpha^{[1]}, \alpha_b^{[1]}, \dots, \alpha^{[L]}, \alpha_b^{[L]}\} \in \text{Tuple}_{\{m'_0, \dots, m'_L\}}$ satisfying some **compatibility conditions** (see Condition 1 and Condition 2), we define a **general embedding** $\mathcal{T}_{\mathcal{I}}^{\alpha} : \text{Tuple}_{\{m_0, \dots, m_L\}} \rightarrow \text{Tuple}_{\{m'_0, \dots, m'_L\}}$ as: For any parameters $\theta_{\text{narr}} = (\mathbf{W}_{\text{narr}}^{[1]}, \mathbf{b}_{\text{narr}}^{[1]}, \dots, \mathbf{W}_{\text{narr}}^{[L]}, \mathbf{b}_{\text{narr}}^{[L]}) \in \text{Tuple}_{\{m_0, \dots, m_L\}}$,*

$$\begin{aligned} \mathcal{T}_{\mathcal{I}}^{\alpha}(\theta_{\text{narr}}) := & \left(\alpha^{[1]} \circ \mathbf{W}_{\text{inter}}^{[1]}, \alpha_b^{[1]} + \mathbf{b}_{\text{inter}}^{[1]}, \dots, \right. \\ & \alpha^{[l]} \circ \mathbf{W}_{\text{inter}}^{[l]}, \alpha_b^{[l]} + \mathbf{b}_{\text{inter}}^{[l]}, \dots, \\ & \left. \alpha^{[L]} \circ \mathbf{W}_{\text{inter}}^{[L]}, \alpha_b^{[L]} + \mathbf{b}_{\text{inter}}^{[L]} \right), \end{aligned}$$

where

$$\mathbf{W}_{\text{inter}}^{[l]} := \left[\left(\mathbf{W}_{\text{narr}}^{[l]} \right)_{\mathcal{I}_l(i), \mathcal{I}_{l-1}(j)} \right], \quad \mathbf{b}_{\text{inter}}^{[l]} := \left(\left(\mathbf{b}_{\text{narr}}^{[l]} \right)_{\mathcal{I}_l(k)} \right)$$

for $l \in [L]$ with $i, k \in [m'_l], j \in [m'_{l-1}]$, and \circ is the Hadamard product.

Remark 5.1. Since $\mathcal{I}_l : [m'_l] \rightarrow [m_l] \sqcup \{0\}$ for $l \in [L-1]$, and the components in $\mathbf{W}_{\text{narr}}^{[l]}$ and $\mathbf{b}_{\text{narr}}^{[l]}$ are not defined for zero indices, i.e., no definitions can be found for $(\mathbf{W}_{\text{narr}}^{[l]})_{0j}$, $(\mathbf{W}_{\text{narr}}^{[l]})_{i0}$, $(\mathbf{b}_{\text{narr}}^{[l]})_0$, for any $l \in [L-1]$ with $i \in [m_l]$ and $j \in [m_{l-1}] \sqcup \{0\}$, for convenience of expression, we set $(\mathbf{W}_{\text{narr}}^{[l]})_{0j} = 1$, $(\mathbf{W}_{\text{narr}}^{[l]})_{i0} = 1$, and $(\mathbf{b}_{\text{narr}}^{[l]})_0 = 0$, with $i \in [m_l]$ and $j \in [m_{l-1}] \sqcup \{0\}$.

Now we proceed to state out the **certain conditions** for the tuple

$$\alpha = \{\alpha^{[1]}, \alpha_b^{[1]}, \dots, \alpha^{[L]}, \alpha_b^{[L]}\} \in \text{Tuple}_{\{m'_0, \dots, m'_L\}} \quad (5.2)$$

in Definition 5.2.

Condition 1 (Compatibility conditions I) (see Fig. 5.1 for illustration). The elements $\{\alpha^{[l]}\}_{l=1}^L$ in (5.2) satisfy that: there exists a collection of auxiliary variables

$$\beta := \left\{ \beta_j^{[l]} \in \mathbb{R} \mid l \in [0 : L], j \in [m'_l] \setminus \mathcal{I}_l^{-1}(0) \right\}$$

such that

- $\beta_k^{[L]} = 1$ for $k \in [m'_L]$. (Since \mathcal{I}_L is the identity map, $\mathcal{I}_L^{-1}(0) = \emptyset$).
- **Forward conditions:**

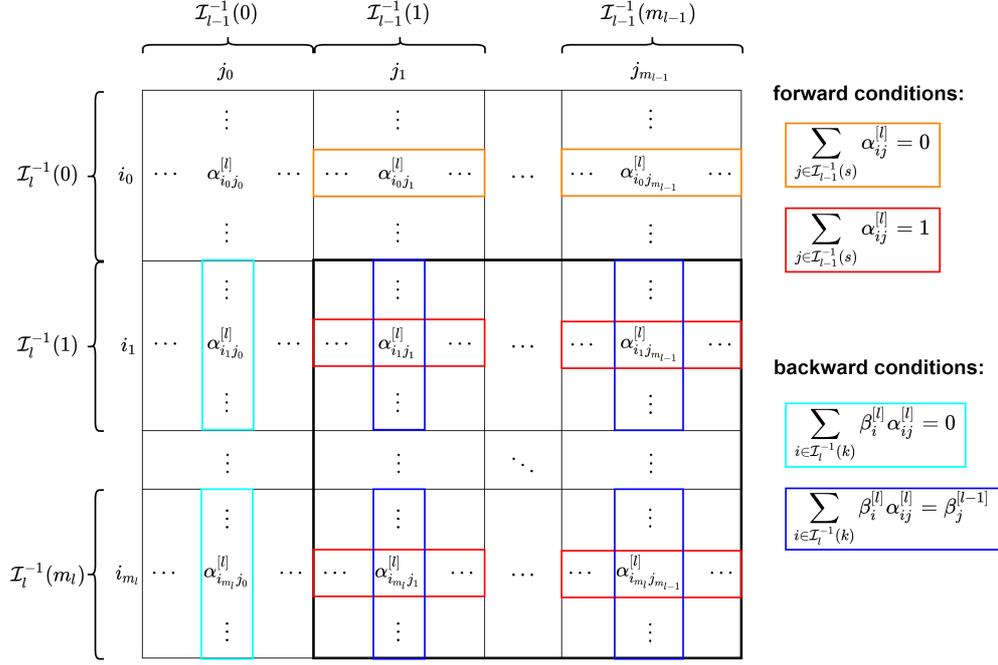


Figure 5.1: Illustration of forward and backward compatibility conditions I for $\alpha^{[l]}$.

- **Effective neurons** $\mathcal{I}_{l-1}^{-1}(s)$ forward to an effective neuron i : For $i \notin \mathcal{I}_l^{-1}(0)$, $s \in [m_{l-1}]$, we have $\sum_{j \in \mathcal{I}_{l-1}^{-1}(s)} \alpha_{ij}^{[l]} = 1$;
- **Effective neurons** $\mathcal{I}_{l-1}^{-1}(s)$ forward to a null neuron i : For $i \in \mathcal{I}_l^{-1}(0)$, $s \in [m_{l-1}]$, we have $\sum_{j \in \mathcal{I}_{l-1}^{-1}(s)} \alpha_{ij}^{[l]} = 0$.

• **Backward conditions:**

- **Effective neurons** $\mathcal{I}_l^{-1}(k)$ backpropagate to an effective neuron j : For $j \notin \mathcal{I}_{l-1}^{-1}(0)$, $k \in [m_l]$, we have $\sum_{i \in \mathcal{I}_l^{-1}(k)} \beta_i^{[l]} \alpha_{ij}^{[l]} = \beta_j^{[l-1]}$;
- **Effective neurons** $\mathcal{I}_l^{-1}(k)$ backpropagate to a null neuron j : For $j \in \mathcal{I}_{l-1}^{-1}(0)$, $k \in [m_l]$, we have $\sum_{i \in \mathcal{I}_l^{-1}(k)} \beta_i^{[l]} \alpha_{ij}^{[l]} = 0$.

The compatibility conditions for $\{\alpha_b^{[l]}\}_{l=1}^L$ is stated in the following. In order for that, we need another collection of auxiliary variables $\mathbf{B}_* := \{(\mathbf{b}_*)_i \in \mathbb{R} \mid l \in [0 : L], i \in \mathcal{I}_l^{-1}(0)\}$. We term \mathbf{B}_* the *effective biases for null neurons*.

Condition 2 (Compatibility conditions II) (see Fig. 5.2 for illustration). The rest of the elements in α , i.e., $\{\alpha_b^{[l]}\}_{l=1}^L$ satisfy that: There also exists a collection of auxiliary variables,

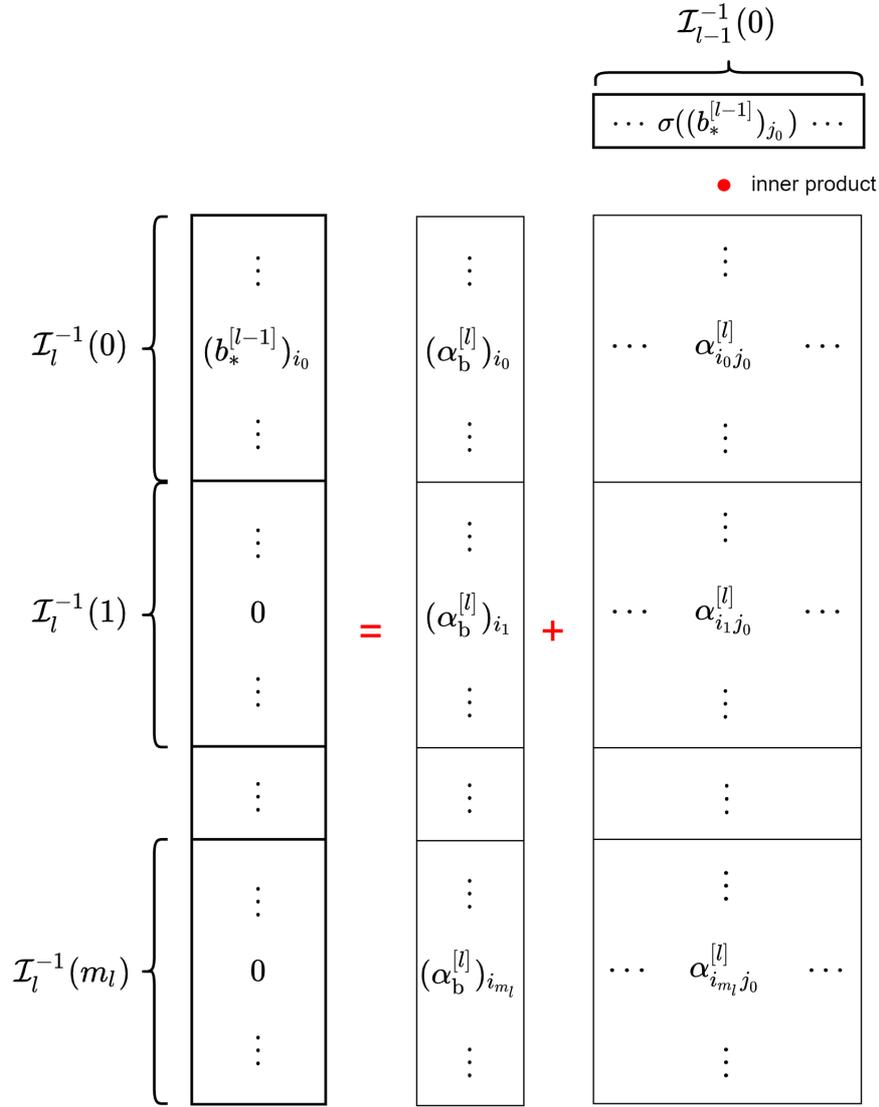


Figure 5.2: Illustration of compatibility conditions II for $\alpha_{\mathbf{b}}^{[l]}$.

termed the *effective biases of null neurons*

$$\mathbf{B}_* := \left\{ (\mathbf{b}_*^{[l]})_i \in \mathbb{R} \mid l \in [0 : L], i \in \mathcal{I}_l^{-1}(0) \right\},$$

we have

- **Effective neurons:** $(\alpha_{\mathbf{b}}^{[l]})_i = 0$ for any $l \in [L]$ with $i \notin \mathcal{I}_l^{-1}(0)$;
- **Null neurons $\mathcal{I}_{l-1}^{-1}(0)$ forward to a null neuron i :** for $i \in \mathcal{I}_l^{-1}(0)$, $s = 0$, we have

Forward

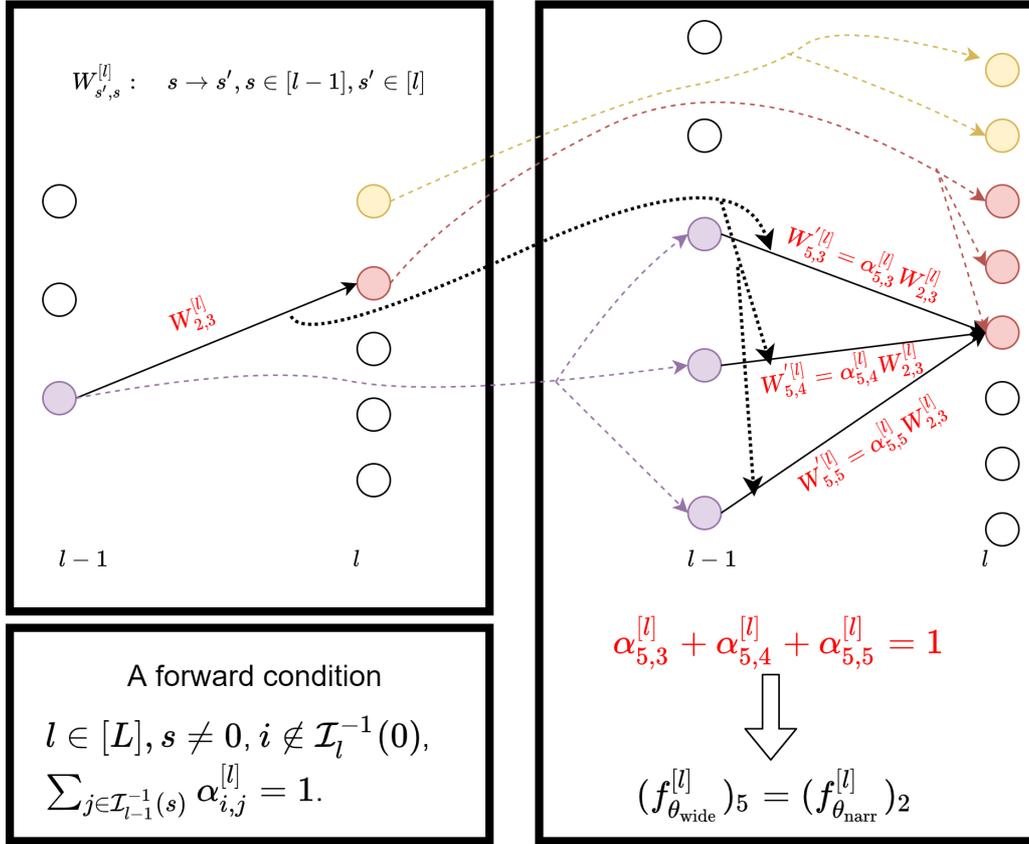


Figure 5.3: Illustration of the forward conditions for split neurons in compatibility conditions I. The sum of the input weights from neurons that are split from the same previous neuron to a post neuron should be equal to the weight from the previous to the post neuron.

$$\sum_{j \in \mathcal{I}_{l-1}^{-1}(0)} \alpha_{ij}^{[l]} \sigma((\mathbf{b}_*^{[l-1]})_j) + (\alpha_{\mathbf{b}}^{[l]})_i = (\mathbf{b}_*^{[l]})_i;$$

- **Null neurons** $\mathcal{I}_{l-1}^{-1}(0)$ **forward** to an effective neuron i : for $i \notin \mathcal{I}_l^{-1}(0)$, $s = 0$, we have $\sum_{j \in \mathcal{I}_{l-1}^{-1}(0)} \alpha_{ij}^{[l]} \sigma((\mathbf{b}_*^{[l-1]})_j) + (\alpha_{\mathbf{b}}^{[l]})_i = 0$.

We then illustrate the general compatible embedding as follows. First, we illustrate an example of the general compatible embedding without null neurons through the forward conditions in Fig. 5.3 and the backward conditions in Fig. 5.4. Second, we illustrate the conditions related to null neurons in Fig. 5.5.

Backward

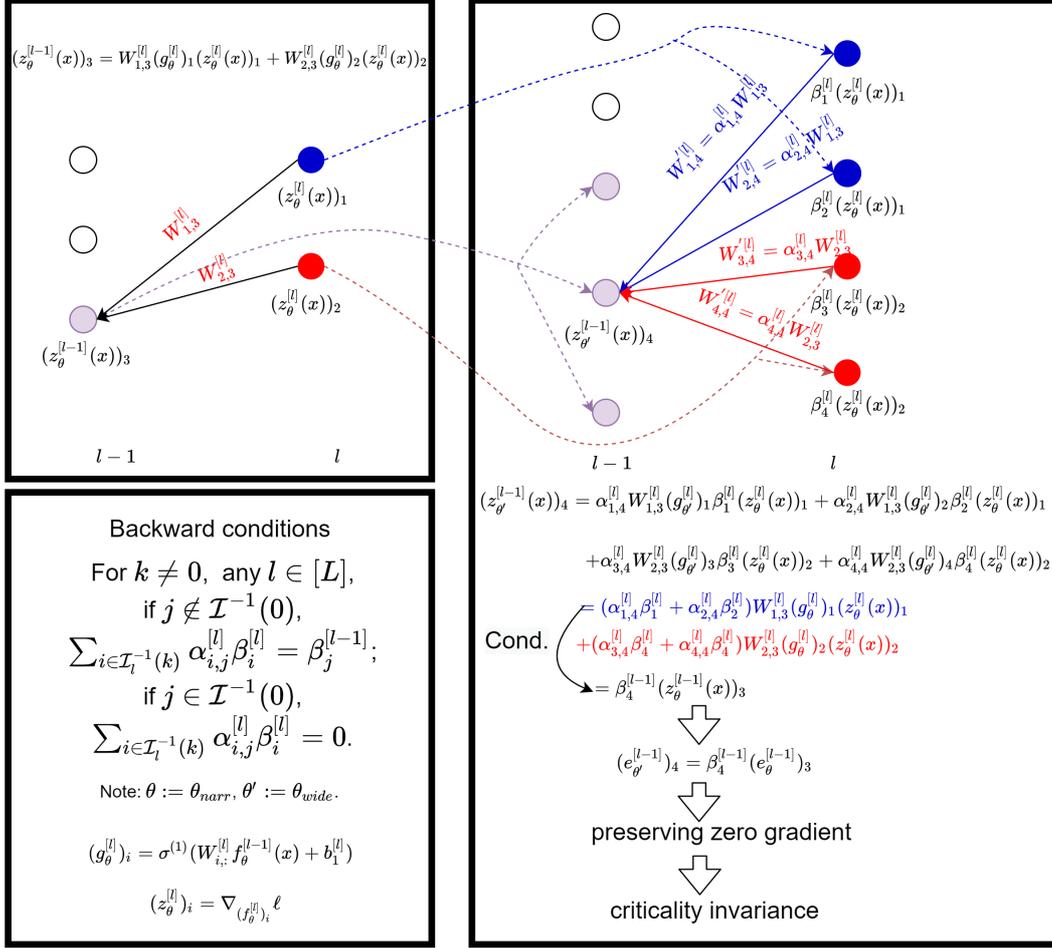


Figure 5.4: Illustration of the backward conditions for split neurons in compatibility conditions I. The gradient of the loss w.r.t. to the neuron output of the wide network is proportional to the gradient of the loss w.r.t. the corresponding neuron in the narrow network.

Theorem 5.1. *General compatible embedding is a critical embedding.*

Remark that we later name it as general compatible critical embedding in this work.

5.1 Special cases of general compatible critical embedding

In the following, we present some special cases of general compatible critical embedding. Specifically, the three-fold global splitting embedding is the key to the proof of a necessary condition of a “truly-bad” critical point in Section 6.

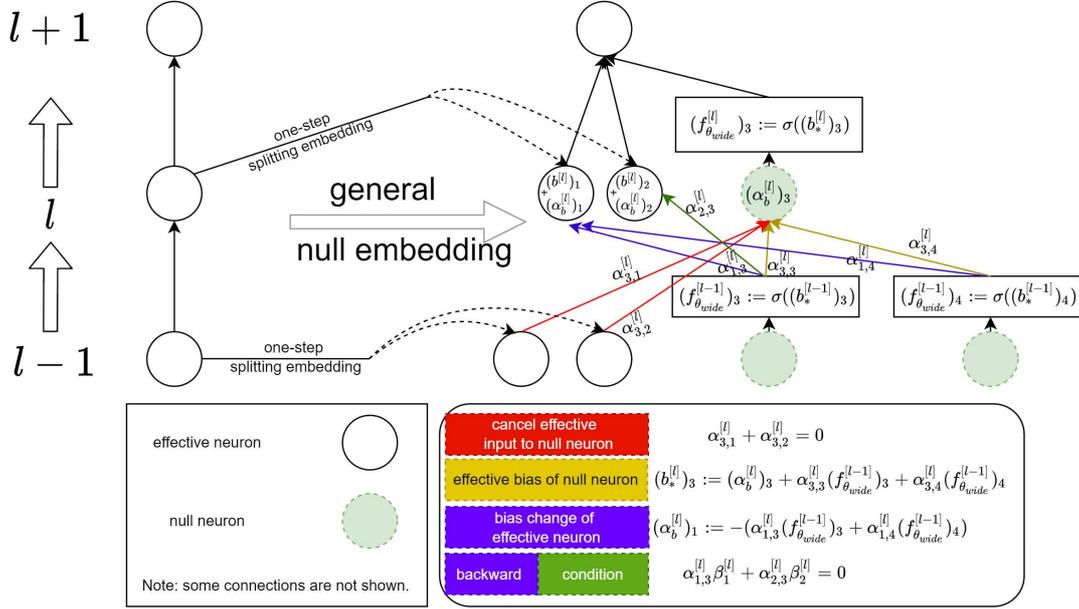


Figure 5.5: Illustration of compatibility conditions for null neurons. Red: The input from effective neurons split from the same neuron should cancel out. Yellow: the effective bias of a null neuron is the sum of its bias (can be any value) and its input from other null neurons. Purple: The added null neurons have constant input to the effective neuron, which would be added a bias correction to cancel the constant input. Purple+Green: the backward condition of a null neuron leads to the zero gradient.

Definition 5.3 (Splitting embedding). For a general compatible critical embedding $\mathcal{T}_{\mathcal{I}}^{\alpha}$, if there is no null neuron in the embedded NN, i.e., $\mathcal{I}_l^{-1}(0) = \emptyset$ for any $l \in [L - 1]$, then we call $\mathcal{T}_{\mathcal{I}}^{\alpha}$ a splitting embedding.

Definition 5.4 (Null embedding). For a general compatible critical embedding $\mathcal{T}_{\mathcal{I}}^{\alpha}$, if only null neurons are added, i.e., $\#\mathcal{I}_l^{-1}(0) = m'_l - m_l$ for any $l \in [L - 1]$, hence there is no neuron splitting, i.e., $\#\mathcal{I}_l^{-1}(s) = 1$ for any $l \in [L - 1], s \in [m_l]$, then we call $\mathcal{T}_{\mathcal{I}}^{\alpha}$ a null embedding.

Remark 5.2. One-step null embedding and its multi-step composition are special cases of null embedding. Similarly, one-step splitting embedding and its multi-step composition are special cases of splitting embedding. Multi-step embedding composed by a mixture of one-step null and splitting embedding is a special case of general compatible critical embedding.

Example 5.1 (A three-fold global splitting embedding). We define the operator $\mathcal{T}_{\text{global}}$ applying on θ as follows

$$\mathcal{T}_{\text{global}}(\theta)|_1 = \left(\begin{bmatrix} \mathbf{W}^{[1]} \\ \mathbf{W}^{[1]} \\ \mathbf{W}^{[1]} \end{bmatrix}, \begin{bmatrix} \mathbf{b}^{[1]} \\ \mathbf{b}^{[1]} \\ \mathbf{b}^{[1]} \end{bmatrix} \right),$$

$$\mathcal{T}_{\text{global}}(\boldsymbol{\theta})|_l = \left(\begin{bmatrix} \mathbf{W}^{[l]} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{W}^{[l]} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{W}^{[l]} \end{bmatrix}, \begin{bmatrix} \mathbf{b}^{[l]} \\ \mathbf{b}^{[l]} \\ \mathbf{b}^{[l]} \end{bmatrix} \right), \quad l \in [2 : L - 1],$$

$$\mathcal{T}_{\text{global}}(\boldsymbol{\theta})|_L = \left(\begin{bmatrix} \mathbf{W}^{[L]} & \mathbf{W}^{[L]} & -\mathbf{W}^{[L]} \end{bmatrix}, \mathbf{b}^{[L]} \right).$$

Remark 5.3. For this global splitting embedding, $m'_l = 3m_l$ and $\mathcal{I}_l^{-1}(s) = \{s, s + m_l, s + 2m_l\}$ for any $l \in [L - 1]$ and $s \in [m_l]$. Moreover, for $l = 1$,

$$\boldsymbol{\alpha}^{[1]} = \begin{bmatrix} \mathbf{1}_{m_1 \times m_0} \\ \mathbf{1}_{m_1 \times m_0} \\ \mathbf{1}_{m_1 \times m_0} \end{bmatrix},$$

$$\boldsymbol{\alpha}_{\mathbf{b}}^{[1]} = \mathbf{0}_{3m_1 \times 1},$$

and for $l \in [2 : L - 1]$, we have

$$\boldsymbol{\alpha}^{[l]} = \begin{bmatrix} \mathbf{1}_{m_l \times m_{l-1}} & \mathbf{0}_{m_l \times m_{l-1}} & \mathbf{0}_{m_l \times m_{l-1}} \\ \mathbf{0}_{m_l \times m_{l-1}} & \mathbf{1}_{m_l \times m_{l-1}} & \mathbf{0}_{m_l \times m_{l-1}} \\ \mathbf{0}_{m_l \times m_{l-1}} & \mathbf{0}_{m_l \times m_{l-1}} & \mathbf{1}_{m_l \times m_{l-1}} \end{bmatrix},$$

$$\boldsymbol{\alpha}_{\mathbf{b}}^{[l]} = \mathbf{0}_{3m_l \times 1},$$

and

$$\boldsymbol{\alpha}^{[L]} = \begin{bmatrix} \mathbf{1}_{m_L \times m_{L-1}} & \mathbf{1}_{m_L \times m_{L-1}} & -\mathbf{1}_{m_L \times m_{L-1}} \end{bmatrix},$$

$$\boldsymbol{\alpha}_{\mathbf{b}}^{[L]} = \mathbf{0}_{m_L \times 1}.$$

Moreover, for any $l \in [L - 1]$, $\beta_j^{[l]} = 1$ when $j \in [2m_l]$ and $\beta_j^{[l]} = -1$ when $j \in [2m_l + 1 : 3m_l]$. It is a special case of splitting embedding.

Example 5.2 (A two-fold global null embedding). We define the operator \mathcal{T}_{GN} applying on $\boldsymbol{\theta}$ as follows

$$\mathcal{T}_{\text{GN}}(\boldsymbol{\theta})|_1 = \left(\begin{bmatrix} \mathbf{W}^{[1]} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \mathbf{b}^{[1]} \\ \mathbf{b}_{**}^{[1]} \end{bmatrix} \right),$$

$$\mathcal{T}_{\text{GN}}(\boldsymbol{\theta})|_l = \left(\begin{bmatrix} \mathbf{W}^{[l]} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}, \begin{bmatrix} \mathbf{b}^{[l]} \\ \mathbf{b}_{**}^{[l]} \end{bmatrix} \right),$$

$$\mathcal{T}_{\text{GN}}(\boldsymbol{\theta})|_L = \left(\begin{bmatrix} \mathbf{W}^{[L]} & \mathbf{0} \end{bmatrix}, \mathbf{b}^{[L]} \right),$$

where $\mathbf{b}_{**}^{[l]}$ takes an arbitrary real vector of the same dimension as $\mathbf{b}^{[l]}$ for any $l \in [L - 1]$.

Remark 5.4. For this global null embedding, $m'_l = 2m_l$, $\mathcal{I}_l^{-1}(s) = \{s\}$ and $\mathcal{I}_l^{-1}(0) = [m_l + 1 : 2m_l]$ for any $l \in [L - 1]$ and $s \in [m_l]$. Then for $l = 1$,

$$\boldsymbol{\alpha}^{[1]} = \begin{bmatrix} \mathbf{1}_{m_1 \times m_0} \\ \mathbf{0}_{m_1 \times m_0} \end{bmatrix},$$

$$\alpha_{\mathbf{b}}^{[1]} = \begin{bmatrix} \mathbf{0}_{m_1 \times 1} \\ \mathbf{b}_{**}^{[1]} \end{bmatrix},$$

and for $l \in [2 : L - 1]$, we have

$$\alpha^{[l]} = \begin{bmatrix} \mathbf{1}_{m_l \times m_{l-1}} & \mathbf{0}_{m_l \times m_{l-1}} \\ \mathbf{0}_{m_l \times m_{l-1}} & \mathbf{0}_{m_l \times m_{l-1}} \end{bmatrix},$$

$$\alpha_{\mathbf{b}}^{[l]} = \begin{bmatrix} \mathbf{0}_{m_l \times 1} \\ \mathbf{b}_{**}^{[l]} \end{bmatrix},$$

and

$$\alpha^{[L]} = \begin{bmatrix} \mathbf{1}_{m_L \times m_{L-1}} & \mathbf{0}_{m_L \times m_{L-1}} \end{bmatrix},$$

$$\alpha_{\mathbf{b}}^{[L]} = \mathbf{0}_{m_L \times 1}.$$

Moreover, for any $l \in [L - 1]$, $(\mathbf{b}_{*}^{[l]})_i = (\mathbf{b}_{**}^{[l]})_{i-m_l}$ for $i \in [m_l + 1 : 2m_l]$. It is a special case of null embedding.

6 Analysis of critical points/submanifolds by critical embeddings

6.1 Degeneracy of critical points/submanifolds

A key observation to the critical embeddings proposed above is that it is not unique. Actually, given any NN and a wider NN, there is a class of critical embeddings. Therefore, any critical point can be embedded to a set of critical points in a wider NN. By a K -step composition embedding, because embedding $\mathcal{T}_{l_i, s_i}^{\alpha_i}$ at each step i has one degree of freedom parameterized by α_i given l_i and s_i , one critical point in general can be embedded to a K -dimensional critical affine subspace, which provides a lower bound to the degeneracy of embedded critical points. Precisely, in [24], we prove the following theorem using composition of one-step splitting embeddings.

Theorem 6.1 (Theorem 2 in [24]). *Given an $\text{NN}(\{m_l\}_{l=0}^L)$ and a K -neuron wider $\text{NN}(\{m'_l\}_{l=0}^L)$, then for any critical point $\theta_{\text{narr}}^c = (\mathbf{W}^{[1]}, \mathbf{b}^{[1]}, \dots, \mathbf{W}^{[L]}, \mathbf{b}^{[L]})$ satisfying $\mathbf{W}^{[l]} \neq \mathbf{0}$ for each $l \in [L]$, θ_{narr}^c can be critically embedded to a K -dimensional critical affine subspace $\mathcal{M}_{\text{wide}} = \{\theta_{\text{wide}} + \sum_{i=1}^K \alpha_i \boldsymbol{\mu}_i \mid \alpha_i \in \mathbb{R}\}$ of loss landscape of $\text{NN}(\{m'_l\}_{l=0}^L)$. Here $\theta_{\text{wide}} := (\prod_{l=1}^K \mathcal{T}_{l, s_l})(\theta_{\text{narr}})$ and $\boldsymbol{\mu}_i := \mathcal{T}_{l_K, s_K} \cdots \mathcal{V}_{l_i, s_i} \cdots \mathcal{T}_{l_1, s_1}(\theta_{\text{narr}})$, where $s_l \neq 0$ for all $l \in [K]$. (The definitions of \mathcal{V}_{l, s_l} and \mathcal{T}_{l, s_l} can be found in Definition 4.4).*

In addition, in this work, as we propose above a wider class of general compatible critical embeddings, we attempt to obtain a better estimate of the dimension of critical submanifolds corresponding to any given critical function $f_{\theta_{\text{narr}}^c}$ of a narrow NN. A gross estimate yields the following heuristic argument.

Heuristic argument. Given any total index mapping \mathcal{I} , the degree of freedom of all possible α for general compatible critical embedding $\mathcal{T}_{\mathcal{I}}^{\alpha}$ to a K -neuron wider NN is $K + \sum_{l \in [L]} K_l K_{l-1}$, where $K_l := m'_l - m_l$, and $K := \sum_{l \in [L-1]} K_l$.

Justification. By the definition of general compatible critical embedding, $\{\alpha^{[l]}\}_{l=1}^L$ and auxiliary variables $\beta = \{\beta_j^{[l]} \in \mathbb{R} \mid l \in [L], j \in [m'_l] \setminus \mathcal{I}_l^{-1}(0)\}$ with $\beta_k^{[L]} = \mathbf{1}$ satisfy forward and backward compatibility conditions.

Step 1: We first observe that $\sum_{i \in \mathcal{I}_l^{-1}(s)} \beta_i^{[l]} = 1$ for any $s \in [m_l]$. Thus, β has $K - m_{\text{null}}$ degrees of freedom, where m_{null} is the number of null neurons.

Step 2: Given any such β , we now consider the degrees of freedom for $\{\alpha^{[l]}\}_{l=1}^L$. Note that forward and backward conditions are not independent of one another because

$$\sum_{j \in \mathcal{I}_{l-1}^{-1}(s)} \sum_{i \in \mathcal{I}_l^{-1}(k)} \beta_i^{[l]} \alpha_{ij}^{[l]} = \sum_{i \in \mathcal{I}_l^{-1}(k)} \beta_i^{[l]} = 1 = \sum_{j \in \mathcal{I}_{l-1}^{-1}(s)} \beta_j^{[l-1]}$$

automatically holds for the given β . Therefore, there are $m'_l m_{l-1} + m_l m'_{l-1} - m_l m_{l-1}$ independent linear equations for $m'_l m'_{l-1}$ parameters in $\alpha^{[l]}$ for each layer l , resulting in $\sum_{l \in [L]} K_l K_{l-1}$ degrees of freedom in total.

Step 3: Given any effective biases B_* as auxiliary variables for all null neurons and any $\alpha^{[l]}$ satisfying forward and backward conditions, $(\alpha_b^{[l]})_i$ is uniquely determined for all neurons in the wide NN. Therefore, there are m_{null} degrees of freedom in B_* .

In the end, the degrees of freedom in α is the summation of degrees of freedom in auxiliary variables β , degrees of freedom in $\{\alpha^{[l]}\}_{l=1}^L$ given β and additional degrees of freedom in B_* for all null neurons, which add up to $K + \sum_{l \in [L]} K_l K_{l-1}$ degrees of freedom.

Intuitively, degrees of freedom in critical embedding \mathcal{T} transform into dimensions of critical submanifolds, resulting in a higher estimate of degrees of degeneracy $K + \sum_{l \in [L]} K_l K_{l-1}$ in comparison to K obtained by multi-step composition embedding in [24]. We note that the nonlinear coupling between α and β in general results in non-affine curved critical submanifolds for three-layer or deeper NNs containing the corresponding K -dimensional critical affine subspaces identified in Theorem 6.1. A more rigorous estimate requires careful handling of the nonlinear coupling. We leave it to later works.

6.2 Irreversible transition to strict-saddle point

In this subsection, we look further into the transition between different types of critical points, e.g., local minima, saddle points, through critical embeddings. We are specifically interested in the transition of a critical point to a strict-saddle point due to its good optimization guarantee detailed in [22]. Strict-saddle point is defined as follows.

Definition 6.1 (Strict-saddle). θ is a strict-saddle point of loss landscape $R_S(\cdot)$ if

(i) $\nabla R_S(\theta) = 0$;

(ii) Hessian matrix $H_S(\theta)$ has at least one negative eigenvalue.

Based on above definition, we prove the following irreversibility property for any critical embedding, which guarantees that a strict-saddle point always embeds to strict-saddle points.

Theorem 6.2. *Given an $\text{NN}(\{m_l\}_{l=0}^L)$ and any of its parameters $\theta \in \mathbb{R}^M$, for any critical embedding $\mathcal{T} : \mathbb{R}^M \rightarrow \mathbb{R}^{M'}$ to any wider $\text{NN}(\{m'_l\}_{l=0}^L)$, the number of positive, zero, negative eigenvalues of $\mathbf{H}_S(\mathcal{T}(\theta))$ is no less than the counterparts of $\mathbf{H}_S(\theta)$.*

6.3 “Truly-bad” critical point

The above irreversibility property of any critical embedding provokes the following thought about a conventional bad local minimum: for any bad local minimum of a given NN, if it can become a strict-saddle point in wider NNs, then it should not be a problem as we can simply use a wider NN in practice. However, there is still a “truly-bad” situation in which a critical point may never become a strict-saddle point through any critical embedding. In the following, through proving a stringent necessary condition for such a “truly-bad” critical point defined below, we justify its rarity, hence providing a potential mechanism to understand the easy optimization of wide NNs widely observed in practice.

We denote hereafter that $A \succeq B$ if and only if $A - B$ is a semi-positive definite matrix, and $A \succ B$ if and only if $A - B$ is a strictly positive definite matrix.

Definition 6.2 (“Truly-bad” critical point). *Given any data S , loss $\ell(\cdot, \cdot)$ and activation $\sigma(\cdot)$, for any NN, if there exists a critical point $\theta^c \in \Theta^c$ satisfying that:*

- (i) θ^c is not a strict-saddle point [22, Definition 1];
 - (ii) For any critical embedding \mathcal{T} , $\mathcal{T}(\theta^c)$ is also not a strict-saddle point,
- then we term this critical point a “truly-bad” critical point.

We would like to introduce some additional notations in order to state Lemma 6.1 and Lemma 6.2. We denote

$$\mathbf{H}_S^{(1),[L-1]}(\theta) := \sum_{i,j=1}^{m_L} \mathbb{E}_S \partial_{ij} \ell(\mathbf{f}_\theta, \mathbf{f}^*) \nabla_{\theta^{[L-1]}}(\mathbf{f}_\theta)_i (\nabla_{\theta^{[L-1]}}(\mathbf{f}_\theta)_j)^\top,$$

and

$$\mathbf{H}_S^{(2),[L-1]}(\theta) := \sum_{i,j=1}^{m_L} \mathbb{E}_S \partial_{ij} \ell(\mathbf{f}_\theta, \mathbf{f}^*) \mathbf{W}_{ij}^{[L]} \nabla_{\theta^{[L-1]}} \nabla_{\theta^{[L-1]}} \left(\mathbf{f}_\theta^{[L-1]} \right)_j.$$

We denote further that $\mathbf{H}_S^{[L-1]}(\theta) := \mathbf{H}_S^{(1),[L-1]}(\theta) + \mathbf{H}_S^{(2),[L-1]}(\theta)$. Obviously, $\mathbf{H}_S^{(1),[L-1]}(\theta)$, $\mathbf{H}_S^{(2),[L-1]}(\theta)$, $\mathbf{H}_S^{[L-1]}(\theta) \in \mathbb{R}^{M^{[L-1]} \times M^{[L-1]}}$, and we state Lemma 6.1 and Lemma 6.2 as follows.

Lemma 6.1. *Given any data S , loss $\ell(\cdot, \cdot)$ and activation $\sigma(\cdot)$, for any NN, if a critical point $\theta^c \in \Theta^c$ satisfies:*

- (i) $\mathbf{H}_S(\theta^c) \succeq 0$;

$$(ii) \mathbf{H}_S^{(2),[L-1]}(\boldsymbol{\theta}^c) \neq \mathbf{0},$$

then there exists a general compatible critical embedding \mathcal{T} , such that $\mathcal{T}(\boldsymbol{\theta}^c)$ is a strict-saddle point.

Lemma 6.2. Given any data S , loss $\ell(\cdot, \cdot)$ and activation $\sigma(\cdot)$, for any NN, if a critical point $\boldsymbol{\theta}^c \in \Theta^c$ satisfies:

$$(i) \mathbf{H}_S(\boldsymbol{\theta}^c) \succeq \mathbf{0};$$

$$(ii) \mathbf{H}_S^{(2),[L-1]}(\boldsymbol{\theta}^c) = \mathbf{0};$$

$$(iii) \mathbf{H}_S^{(2)}(\boldsymbol{\theta}^c) \neq \mathbf{0},$$

then there exists a general compatible critical embedding \mathcal{T} , such that $\mathcal{T}(\boldsymbol{\theta}^c)$ is a strict-saddle point.

Theorem 6.3. Given any data S , loss $\ell(\cdot, \cdot)$ and activation $\sigma(\cdot)$, for any NN, if a critical point $\boldsymbol{\theta}^c$ with $\mathbf{H}_S(\boldsymbol{\theta}^c) \succeq \mathbf{0}$ satisfies $\mathbf{H}_S^{(2)}(\boldsymbol{\theta}^c) \neq \mathbf{0}$, then there exists a general compatible critical embedding \mathcal{T} such that $\mathcal{T}(\boldsymbol{\theta}^c)$ is a strict-saddle point, i.e., $\mathbf{H}_S(\mathcal{T}(\boldsymbol{\theta}^c))$ has at least one negative eigenvalue.

Proof. This can be directly obtained by Lemmas 6.1 and 6.2. \square

Theorem (short version of Theorem 6.3). $\mathbf{H}_S^{(2)}(\boldsymbol{\theta}^c) = \mathbf{0}$ is a necessary condition for a critical point $\boldsymbol{\theta}^c$ being a "truly-bad" critical point.

7 Conclusion and discussion

In this work, we prove the Embedding Principle that loss landscape of an NN *contains* all critical points/functions of all the narrower NNs. We define the critical embedding, which serves as the key tool not only to the proof of Embedding Principle but also to the study of the general geometry of loss landscape. Importantly, we discover a wide class of general compatible embedding, by which we obtain rich understanding about the critical points/submanifolds, e.g., lower bound of their degree of degeneracy, their easy and irreversible transition to strict-saddle points.

The general compatible embedding proposed in this work unravels that the critical embedding in general is not limited to the composition of one-step embeddings, but instead, it can be a collective operation. As a consequence, all critical points embedded from a narrower NN form high-dimensional critical submanifolds, which in general are not affine subspaces for three-layer or deeper NNs. It is interesting and important. However, it remains a problem for the future study about whether the general compatible embeddings in certain sense are all the critical embeddings.

Embedding Principle provides an integrated view about loss landscapes of NNs with different widths. Specifically, it informs us that a specific bad local minimum in a NN with a fixed width may not be a big deal for optimization as long as it can become a strict-saddle point in wider NNs, i.e., not a "truly-bad" critical point. In this work, we prove a stringent necessary condition for a "truly-bad" critical point. Still, it remains a problem

about whether non-trivial “truly-bad” critical points indeed exist, and whether they are indeed a headache for optimization. We remark that complementary to the study of the “truly-bad” critical points, we also need to further study whether NNs can generate new bad local-minima forever when they become wider and wider to better understand the easy optimization of wide NNs.

We emphasize that Embedding Principle provides a function space view on the critical points /submanifolds of loss landscape, which is of great importance for studying the implicit regularization and generalization of NNs. In recent years, it has been more and more clear that optimization and generalization are heavily intertwined with one another for deep learning. Therefore, one can not expect to develop a deep learning theory with optimization theory and generalization theory established separately. Yet, for a long time, the study of loss landscape focuses mainly on the parameter space in pursuit of an optimization guarantee. On the other hand, the study of generalization focuses mainly on the function space, failing to incorporate the geometry of loss landscape in parameter space which is key to the training. Now, by the Embedding Principle, we see a hierarchical structure of critical functions of different complexities of the loss landscape, which originate from loss landscape with clear optimization implication while being amenable to the analysis of complexity with clear generalization implication. Such a critical function hierarchy reflects the degree of matching between the NN architecture and data, i.e., if critical functions of narrow NNs attain low training error, then the NN architecture well matches the target function and may obtain a well generalized solution through training like in Fig. 4.3. Even though more works need to be done to unravel the full details and implication of this critical function hierarchy, we believe it is an important piece and may be a key to the deep learning theory.

Acknowledgments

This work is sponsored by the National Key R&D Program of China Grant No. 2019YFA0709503 (Z. X.), the Shanghai Sailing Program, the Natural Science Foundation of Shanghai Grant No. 20ZR1429000 (Z. X.), the National Natural Science Foundation of China Grant No. 62002221 (Z. X.), the National Natural Science Foundation of China Grant No. 12101401 (T. L.), the National Natural Science Foundation of China Grant No. 12101402 (Y. Z.), Shanghai Municipal of Science and Technology Project Grant No. 20JC1419500 (Y.Z.), Shanghai Municipal of Science and Technology Major Project No. 2021SHZDZX0102, and the HPC of School of Mathematical Sciences and the Student Innovation Center at Shanghai Jiao Tong University.

A Theoretical details

Lemma (Lemma 4.1 in main text.). For any one-step null embedding $\mathcal{T}_{l,0}^\alpha$, given any $\text{NN}(\{m_l\}_{l=0}^L)$ and its parameters $\theta_{\text{narr}} \in \text{Tuple}_{\{m_0, \dots, m_L\}}$ with $\text{Tuple}_{\{m_0, \dots, m_L\}} \in \mathcal{D}_{l,0}$, we have $\theta_{\text{wide}} := \mathcal{T}_{l,0}^\alpha(\theta_{\text{narr}})$ satisfies the following conditions: given any data S , loss $\ell(\cdot, \cdot)$

and activation $\sigma(\cdot)$, for any $l \in [L - 1]$,

(i) feature vectors in

$$\mathbf{F}_{\boldsymbol{\theta}_{\text{wide}}}: \mathbf{f}_{\boldsymbol{\theta}_{\text{wide}}}^{[l']} = \mathbf{f}_{\boldsymbol{\theta}_{\text{narr}}}^{[l']}, \text{ for } l' \in [L] \text{ and } l' \neq l, \mathbf{f}_{\boldsymbol{\theta}_{\text{wide}}}^{[l]} = \left[(\mathbf{f}_{\boldsymbol{\theta}_{\text{narr}}}^{[l]})^\top, \sigma(\alpha) \right]^\top;$$

(ii) feature gradients in

$$\mathbf{G}_{\boldsymbol{\theta}_{\text{wide}}}: \mathbf{g}_{\boldsymbol{\theta}_{\text{wide}}}^{[l']} = \mathbf{g}_{\boldsymbol{\theta}_{\text{narr}}}^{[l']}, \text{ for } l' \in [L] \text{ and } l' \neq l, \mathbf{g}_{\boldsymbol{\theta}_{\text{wide}}}^{[l]} = \left[(\mathbf{g}_{\boldsymbol{\theta}_{\text{narr}}}^{[l]})^\top, \sigma^{(1)}(\alpha) \right]^\top;$$

(iii) error vectors in

$$\mathbf{Z}_{\boldsymbol{\theta}_{\text{wide}}}: \mathbf{z}_{\boldsymbol{\theta}_{\text{wide}}}^{[l']} = \mathbf{z}_{\boldsymbol{\theta}_{\text{narr}}}^{[l']}, \text{ for } l' \in [L] \text{ and } l' \neq l, \mathbf{z}_{\boldsymbol{\theta}_{\text{wide}}}^{[l]} = \left[\mathbf{z}_{\boldsymbol{\theta}_{\text{narr}}}^{[l]}, 0 \right]^\top;$$

(iv) $\mathcal{T}_{l,0}^\alpha$ is injective for all α ;

(v) $\mathcal{T}_{l,0}^\alpha$ is an affine embedding for all α .

Proof. (i) By the construction of $\boldsymbol{\theta}_{\text{wide}}$, it is clear that $\mathbf{f}_{\boldsymbol{\theta}_{\text{wide}}}^{[l']} = \mathbf{f}_{\boldsymbol{\theta}_{\text{narr}}}^{[l']}$ for all $l' \in [l - 1]$. Then

$$\mathbf{f}_{\boldsymbol{\theta}_{\text{wide}}}^{[l]} = \sigma \left(\begin{bmatrix} \mathbf{W}^{[l]} \\ \mathbf{0}_{1 \times m_{l-1}} \end{bmatrix} \mathbf{f}_{\boldsymbol{\theta}_{\text{narr}}}^{[l-1]} + \begin{bmatrix} \mathbf{b}^{[l]} \\ \alpha \end{bmatrix} \right) = \begin{bmatrix} \mathbf{f}_{\boldsymbol{\theta}_{\text{narr}}}^{[l]} \\ \sigma(\alpha) \end{bmatrix}. \quad (\text{A.1})$$

Note that since

$$\alpha \begin{bmatrix} \mathbf{0}_{m_{l+1} \times (m_l + 1)} \end{bmatrix} \begin{bmatrix} \mathbf{f}_{\boldsymbol{\theta}_{\text{narr}}}^{[l]} \\ \sigma(\alpha) \end{bmatrix} = \mathbf{0}_{m_{l+1} \times 1}.$$

Thus

$$\mathbf{f}_{\boldsymbol{\theta}_{\text{wide}}}^{[l+1]} = \sigma \left(\begin{bmatrix} \mathbf{W}^{[l+1]}, \mathbf{0}_{m_{l+1} \times 1} \end{bmatrix} \begin{bmatrix} \mathbf{f}_{\boldsymbol{\theta}_{\text{narr}}}^{[l]} \\ \sigma(\alpha) \end{bmatrix} + \mathbf{0}_{m_{l+1} \times 1} + \mathbf{b}^{[l+1]} + \alpha \mathbf{0}_{m_{l+1} \times 1} \right) = \mathbf{f}_{\boldsymbol{\theta}_{\text{narr}}}^{[l+1]}. \quad (\text{A.2})$$

Next, by the construction of $\boldsymbol{\theta}_{\text{wide}}$, it is clear that $\mathbf{f}_{\boldsymbol{\theta}_{\text{wide}}}^{[l']} = \mathbf{f}_{\boldsymbol{\theta}_{\text{narr}}}^{[l']}$ for any $l' \in [l + 1 : L]$.

(ii) The results for feature gradients $\mathbf{g}_{\boldsymbol{\theta}_{\text{wide}}}^{[l']} = \mathbf{g}_{\boldsymbol{\theta}_{\text{narr}}}^{[l']}$ for $l' \in [L]$ can be calculated in a similar way except by replacing $\sigma(\cdot)$ with $\sigma^{(1)}(\cdot)$.

(iii) By the backpropagation and the above facts in (i), we have

$$\mathbf{z}_{\boldsymbol{\theta}_{\text{wide}}}^{[L]} = \nabla \ell(\mathbf{f}_{\boldsymbol{\theta}_{\text{wide}}}^{[L]}, \mathbf{y}) = \nabla \ell(\mathbf{f}_{\boldsymbol{\theta}_{\text{narr}}}^{[L]}, \mathbf{y}) = \mathbf{z}_{\boldsymbol{\theta}_{\text{narr}}}^{[L]}.$$

Recall the recurrence relation for $l' \in [l + 1 : L - 1]$, then we recursively obtain the following equality for l' from $L - 1$ down to $l + 1$:

$$\mathbf{z}_{\boldsymbol{\theta}_{\text{wide}}}^{[l']} = (\mathbf{W}^{[l'+1]})^\top \left(\mathbf{z}_{\boldsymbol{\theta}_{\text{wide}}}^{[l'+1]} \circ \mathbf{g}_{\boldsymbol{\theta}_{\text{wide}}}^{[l'+1]} \right) = (\mathbf{W}^{[l'+1]})^\top \left(\mathbf{z}_{\boldsymbol{\theta}_{\text{narr}}}^{[l'+1]} \circ \mathbf{g}_{\boldsymbol{\theta}_{\text{narr}}}^{[l'+1]} \right) = \mathbf{z}_{\boldsymbol{\theta}_{\text{narr}}}^{[l']}. \quad (\text{A.3})$$

Next,

$$\begin{aligned} \mathbf{z}_{\boldsymbol{\theta}_{\text{wide}}}^{[l]} &= \left(\left[\mathbf{W}^{[l+1]}, \mathbf{0}_{m_{l+1} \times 1} \right] + \alpha \left[\mathbf{0}_{m_{l+1} \times (m_l+1)} \right] \right)^\top \left(\mathbf{z}_{\boldsymbol{\theta}_{\text{wide}}}^{[l+1]} \circ \mathbf{g}_{\boldsymbol{\theta}_{\text{wide}}}^{[l+1]} \right) \\ &= \begin{bmatrix} \mathbf{z}_{\boldsymbol{\theta}_{\text{narr}}}^{[l]} \\ 0 \end{bmatrix} + \left[\mathbf{0}_{(m_l+1) \times 1} \right] = \begin{bmatrix} \mathbf{z}_{\boldsymbol{\theta}_{\text{narr}}}^{[l]} \\ 0 \end{bmatrix}^\top. \end{aligned} \quad (\text{A.4})$$

Finally,

$$\begin{aligned} \mathbf{z}_{\boldsymbol{\theta}_{\text{wide}}}^{[l-1]} &= \left[\mathbf{W}^{[l]^\top}, \mathbf{0}_{m_{l-1} \times 1} \right] \left(\begin{bmatrix} \mathbf{z}_{\boldsymbol{\theta}_{\text{narr}}}^{[l]} \\ 0 \end{bmatrix} \circ \begin{bmatrix} \mathbf{g}_{\boldsymbol{\theta}_{\text{narr}}}^{[l]} \\ \sigma^{(1)}(\alpha) \end{bmatrix} \right) \\ &= (\mathbf{W}^{[l]})^\top \left(\mathbf{z}_{\boldsymbol{\theta}_{\text{narr}}}^{[l]} \circ \mathbf{g}_{\boldsymbol{\theta}_{\text{narr}}}^{[l]} \right) + \mathbf{0}_{m_{l-1} \times 1} = \mathbf{z}_{\boldsymbol{\theta}_{\text{narr}}}^{[l-1]}. \end{aligned} \quad (\text{A.5})$$

This with the recurrence relation once again leads to $\mathbf{z}_{\boldsymbol{\theta}_{\text{wide}}}^{[l']} = \mathbf{z}_{\boldsymbol{\theta}_{\text{narr}}}^{[l']}$ for all $l' \in [l-1]$.

(iv) If for $\boldsymbol{\theta}_1, \boldsymbol{\theta}_2 \in \text{Tuple}_{\{m_0, \dots, m_L\}}$ and $\boldsymbol{\theta}_1 \neq \boldsymbol{\theta}_2$, since $\mathcal{T}_{l,0}^\alpha|_k$ for $k \neq l, l+1$ is the identity map, then if there exists some $k_0 \neq l, l+1$, such that $\mathbf{W}_1^{[k_0]} \neq \mathbf{W}_2^{[k_0]}$ or $\mathbf{b}_1^{[k_0]} \neq \mathbf{b}_2^{[k_0]}$, then obviously $\mathcal{T}_{l,0}^\alpha(\boldsymbol{\theta}_1) \neq \mathcal{T}_{l,0}^\alpha(\boldsymbol{\theta}_2)$. If $k_0 = l$ or $k_0 = l+1$, by similar reasoning, $\mathcal{T}_{l,0}^\alpha(\boldsymbol{\theta}_1) \neq \mathcal{T}_{l,0}^\alpha(\boldsymbol{\theta}_2)$.

(v) For $\boldsymbol{\theta}_0 = (\mathbf{W}_0^{[1]}, \mathbf{b}_0^{[1]}, \dots, \mathbf{W}_0^{[L]}, \mathbf{b}_0^{[L]}) \in \text{Tuple}_{\{m_0, \dots, m_L\}}$, we have

$$\begin{aligned} &\tilde{\mathcal{T}}_{l,0}^\alpha(\boldsymbol{\theta}_0) \\ &:= \mathcal{T}_{l,0}^\alpha(\boldsymbol{\theta}_0) - \mathcal{T}_{l,0}^\alpha(\mathbf{0}) \\ &= \left(\mathbf{W}_0^{[1]}, \mathbf{b}_0^{[1]}, \dots, \begin{bmatrix} \mathbf{W}_0^{[l]} \\ \mathbf{0}_{1 \times m_{l-1}} \end{bmatrix}, \begin{bmatrix} \mathbf{b}_0^{[l]} \\ 0 \end{bmatrix}, \begin{bmatrix} \mathbf{W}_0^{[l+1]}, \mathbf{0}_{m_{l+1} \times 1} \end{bmatrix}, \mathbf{b}_0^{[l+1]}, \dots, \mathbf{W}_0^{[L]}, \mathbf{b}_0^{[L]} \right), \end{aligned}$$

obviously $\tilde{\mathcal{T}}_{l,0}^\alpha$ is a linear operator, thus $\mathcal{T}_{l,0}^\alpha$ is an affine operator. \square

Lemma (Lemma 4.2 in main text). For any one-step splitting embedding $\mathcal{T}_{l,s}^\alpha$, given any $\text{NN}(\{m_l\}_{l=0}^L)$ and its parameters $\boldsymbol{\theta}_{\text{narr}} \in \text{Tuple}_{\{m_0, \dots, m_L\}}$ with $\text{Tuple}_{\{m_0, \dots, m_L\}} \in \mathcal{D}_{l,s}$, we have $\boldsymbol{\theta}_{\text{wide}} := \mathcal{T}_{l,s}^\alpha(\boldsymbol{\theta}_{\text{narr}})$ satisfies the following conditions: given any data S , loss $\ell(\cdot, \cdot)$ and activation $\sigma(\cdot)$, for any $l \in [L-1]$,

(i) feature vectors in

$$\mathbf{F}_{\boldsymbol{\theta}_{\text{wide}}}: \mathbf{f}_{\boldsymbol{\theta}_{\text{wide}}}^{[l']} = \mathbf{f}_{\boldsymbol{\theta}_{\text{narr}}}^{[l']}, \text{ for } l' \in [L] \text{ and } l' \neq l, \mathbf{f}_{\boldsymbol{\theta}_{\text{wide}}}^{[l]} = \left[(\mathbf{f}_{\boldsymbol{\theta}_{\text{narr}}}^{[l]})^\top, (\mathbf{f}_{\boldsymbol{\theta}_{\text{narr}}}^{[l]})_s \right]^\top;$$

(ii) feature gradients in

$$\mathbf{G}_{\boldsymbol{\theta}_{\text{wide}}}: \mathbf{g}_{\boldsymbol{\theta}_{\text{wide}}}^{[l']} = \mathbf{g}_{\boldsymbol{\theta}_{\text{narr}}}^{[l']}, \text{ for } l' \in [L] \text{ and } l' \neq l, \mathbf{g}_{\boldsymbol{\theta}_{\text{wide}}}^{[l]} = \left[(\mathbf{g}_{\boldsymbol{\theta}_{\text{narr}}}^{[l]})^\top, (\mathbf{g}_{\boldsymbol{\theta}_{\text{narr}}}^{[l]})_s \right]^\top;$$

(iii) error vectors in $\mathbf{Z}_{\theta_{\text{wide}}}$: $\mathbf{z}_{\theta_{\text{wide}}}^{[l']} = \mathbf{z}_{\theta_{\text{narr}}}^{[l']}$, for $l' \in [L]$ and $l' \neq l$,

$$\mathbf{z}_{\theta_{\text{wide}}}^{[l]} = \left[\left(\mathbf{z}_{\theta_{\text{narr}}}^{[l]} \right)_{[1:s-1]}^\top, (1-\alpha)(\mathbf{z}_{\theta_{\text{narr}}}^{[l]})_s, \left(\mathbf{z}_{\theta_{\text{narr}}}^{[l]} \right)_{[s+1:m_l]}^\top, \alpha(\mathbf{z}_{\theta_{\text{narr}}}^{[l]})_s \right]^\top;$$

(iv) $\mathcal{T}_{l,s}^\alpha$ is injective for all α ;

(v) $\mathcal{T}_{l,s}^\alpha$ is an affine embedding for all α .

Proof. (i) By the construction of θ_{wide} , it is clear that $\mathbf{f}_{\theta_{\text{wide}}}^{[l']} = \mathbf{f}_{\theta_{\text{narr}}}^{[l']}$ for all $l' \in [l-1]$. Then

$$\mathbf{f}_{\theta_{\text{wide}}}^{[l]} = \sigma \left(\left[\begin{array}{c} \mathbf{W}^{[l]} \\ \mathbf{W}_{s,[1:m_{l-1}]}^{[l]} \end{array} \right] \mathbf{f}_{\theta_{\text{narr}}}^{[l-1]} + \left[\begin{array}{c} \mathbf{b}^{[l]} \\ \mathbf{b}_s^{[l]} \end{array} \right] \right) = \left[\begin{array}{c} \mathbf{f}_{\theta_{\text{narr}}}^{[l]} \\ (\mathbf{f}_{\theta_{\text{narr}}}^{[l]})_s \end{array} \right]. \quad (\text{A.6})$$

Note that

$$\alpha \left[\mathbf{0}_{m_{l+1} \times (s-1)}, -\mathbf{W}_{[1:m_{l+1}],s'}^{[l+1]}, \mathbf{0}_{m_{l+1} \times (m_l - s)}, \mathbf{W}_{[1:m_{l+1}],s}^{[l+1]} \right] \left[\begin{array}{c} \mathbf{f}_{\theta_{\text{narr}}}^{[l]} \\ (\mathbf{f}_{\theta_{\text{narr}}}^{[l]})_s \end{array} \right] = \mathbf{0}_{m_{l+1} \times 1}.$$

Thus

$$\mathbf{f}_{\theta_{\text{wide}}}^{[l+1]} = \sigma \left(\left[\mathbf{W}^{[l+1]}, \mathbf{0}_{m_{l+1} \times 1} \right] \left[\begin{array}{c} \mathbf{f}_{\theta_{\text{narr}}}^{[l]} \\ (\mathbf{f}_{\theta_{\text{narr}}}^{[l]})_s \end{array} \right] + \mathbf{0}_{m_{l+1} \times 1} + \mathbf{b}^{[l+1]} + \alpha \mathbf{0}_{m_{l+1} \times 1} \right) = \mathbf{f}_{\theta_{\text{narr}}}^{[l+1]}. \quad (\text{A.7})$$

Next, by the construction of θ_{wide} , it is clear that $\mathbf{f}_{\theta_{\text{wide}}}^{[l']} = \mathbf{f}_{\theta_{\text{narr}}}^{[l']}$ for any $l' \in [l+1:L]$.

(ii) The results for feature gradients $\mathbf{g}_{\theta_{\text{wide}}}^{[l']} = \mathbf{g}_{\theta_{\text{narr}}}^{[l']}$ for $l' \in [L]$ can be calculated in a similar way except by replacing $\sigma(\cdot)$ with $\sigma^{(1)}(\cdot)$.

(iii) By the backpropagation and the above facts in (i), we have

$$\mathbf{z}_{\theta_{\text{wide}}}^{[L]} = \nabla \ell(\mathbf{f}_{\theta_{\text{wide}}}^{[L]}, \mathbf{y}) = \nabla \ell(\mathbf{f}_{\theta_{\text{narr}}}^{[L]}, \mathbf{y}) = \mathbf{z}_{\theta_{\text{narr}}}^{[L]}.$$

Recall the recurrence relation for $l' \in [l+1:L-1]$, then we recursively obtain the following equality for l' from $L-1$ down to $l+1$:

$$\mathbf{z}_{\theta_{\text{wide}}}^{[l']} = (\mathbf{W}^{[l'+1]})^\top \left(\mathbf{z}_{\theta_{\text{wide}}}^{[l'+1]} \circ \mathbf{g}_{\theta_{\text{wide}}}^{[l'+1]} \right) = (\mathbf{W}^{[l'+1]})^\top \left(\mathbf{z}_{\theta_{\text{narr}}}^{[l'+1]} \circ \mathbf{g}_{\theta_{\text{narr}}}^{[l'+1]} \right) = \mathbf{z}_{\theta_{\text{narr}}}^{[l']}. \quad (\text{A.8})$$

Next,

$$\mathbf{z}_{\theta_{\text{wide}}}^{[l]} = \left(\left[\mathbf{W}^{[l+1]}, \mathbf{0}_{m_{l+1} \times 1} \right] + \alpha \left[\mathbf{0}_{m_{l+1} \times (s-1)}, -\mathbf{W}_{[1:m_{l+1}],s'}^{[l+1]}, \mathbf{0}_{m_{l+1} \times (m_l - s)}, \mathbf{W}_{[1:m_{l+1}],s}^{[l+1]} \right] \right)^\top \left(\mathbf{z}_{\theta_{\text{wide}}}^{[l+1]} \circ \mathbf{g}_{\theta_{\text{wide}}}^{[l+1]} \right)$$

$$\begin{aligned}
&= \begin{bmatrix} \mathbf{z}_{\boldsymbol{\theta}_{\text{narr}}}^{[l]} \\ \mathbf{0} \end{bmatrix} + \begin{bmatrix} \mathbf{0}_{(s-1) \times 1} \\ -\alpha(\mathbf{z}_{\boldsymbol{\theta}_{\text{narr}}}^{[l]})_s \\ \mathbf{0}_{(m_l-s) \times 1} \\ \alpha(\mathbf{z}_{\boldsymbol{\theta}_{\text{narr}}}^{[l]})_s \end{bmatrix} \\
&= \left[(\mathbf{z}_{\boldsymbol{\theta}_{\text{narr}}}^{[l]})_{[1:s-1]}^\top, (1-\alpha)(\mathbf{z}_{\boldsymbol{\theta}_{\text{narr}}}^{[l]})_s, (\mathbf{z}_{\boldsymbol{\theta}_{\text{narr}}}^{[l]})_{[s+1:m_l]}^\top, \alpha(\mathbf{z}_{\boldsymbol{\theta}_{\text{narr}}}^{[l]})_s \right]^\top. \tag{A.9}
\end{aligned}$$

Finally,

$$\begin{aligned}
\mathbf{z}_{\boldsymbol{\theta}_{\text{wide}}}^{[l-1]} &= \left[(\mathbf{W}^{[l]})^\top, (\mathbf{W}^{[l]})_{s,[1:m_{l-1}]}^\top \right] \left(\left(\begin{bmatrix} \mathbf{z}_{\boldsymbol{\theta}_{\text{narr}}}^{[l]} \\ \mathbf{0} \end{bmatrix} + \begin{bmatrix} \mathbf{0}_{(s-1) \times 1} \\ -\alpha(\mathbf{z}_{\boldsymbol{\theta}_{\text{narr}}}^{[l]})_s \\ \mathbf{0}_{(m_l-s) \times 1} \\ \alpha(\mathbf{z}_{\boldsymbol{\theta}_{\text{narr}}}^{[l]})_s \end{bmatrix} \right) \circ \begin{bmatrix} \mathbf{g}_{\boldsymbol{\theta}_{\text{narr}}}^{[l]} \\ (\mathbf{g}_{\boldsymbol{\theta}_{\text{narr}}}^{[l]})_s \end{bmatrix} \right) \\
&= (\mathbf{W}^{[l]})^\top \left(\mathbf{z}_{\boldsymbol{\theta}_{\text{narr}}}^{[l]} \circ \mathbf{g}_{\boldsymbol{\theta}_{\text{narr}}}^{[l]} \right) + \mathbf{0}_{m_{l-1} \times 1} = \mathbf{z}_{\boldsymbol{\theta}_{\text{narr}}}^{[l-1]}. \tag{A.10}
\end{aligned}$$

This with the recurrence relation once again leads to $\mathbf{z}_{\boldsymbol{\theta}_{\text{wide}}}^{[l']} = \mathbf{z}_{\boldsymbol{\theta}_{\text{narr}}}^{[l']}$ for all $l' \in [l-1]$.

(iv) If for $\boldsymbol{\theta}_1, \boldsymbol{\theta}_2 \in \text{Tuple}_{\{m_0, \dots, m_L\}}$ and $\boldsymbol{\theta}_1 \neq \boldsymbol{\theta}_2$, since $\mathcal{T}_{l,s}^\alpha|_k$ for $k \neq l, l+1$ is the identity map, then if there exists some $k_0 \neq l, l+1$, such that $\mathbf{W}_1^{[k_0]} \neq \mathbf{W}_2^{[k_0]}$ or $\mathbf{b}_1^{[k_0]} \neq \mathbf{b}_2^{[k_0]}$, then obviously $\mathcal{T}_{l,s}^\alpha(\boldsymbol{\theta}_1) \neq \mathcal{T}_{l,s}^\alpha(\boldsymbol{\theta}_2)$. If $k_0 = l$ or $k_0 = l+1$, by similar reasoning, $\mathcal{T}_{l,s}^\alpha(\boldsymbol{\theta}_1) \neq \mathcal{T}_{l,s}^\alpha(\boldsymbol{\theta}_2)$.

(v) For $\boldsymbol{\theta}_0 = (\mathbf{W}_0^{[1]}, \mathbf{b}_0^{[1]}, \dots, \mathbf{W}_0^{[L]}, \mathbf{b}_0^{[L]}) \in \text{Tuple}_{\{m_0, \dots, m_L\}}$, we have

$$\begin{aligned}
\tilde{\mathcal{T}}_{l,s}^\alpha(\boldsymbol{\theta}_0) &:= \mathcal{T}_{l,s}^\alpha(\boldsymbol{\theta}_0) - \mathcal{T}_{l,s}^\alpha(\mathbf{0}) \\
&= \left(\mathbf{W}_0^{[1]}, \mathbf{b}_0^{[1]}, \dots, \begin{bmatrix} \mathbf{W}_0^{[l]} \\ (\mathbf{W}_0^{[l]})_{s,[1:m_{l-1}]} \end{bmatrix}, \begin{bmatrix} \mathbf{b}_0^{[l]} \\ (\mathbf{b}_0^{[l]})_s \end{bmatrix}, [(\mathbf{W}_0^{[l+1]})_{[1:m_{l+1}], [1:s-1]}, \\
&\quad (1-\alpha)(\mathbf{W}_0^{[l+1]})_{[1:m_{l+1}], s}, (\mathbf{W}_0^{[l+1]})_{[1:m_{l+1}], [s+1:m_l]}, \alpha(\mathbf{W}_0^{[l+1]})_{[1:m_{l+1}], s}, \\
&\quad \mathbf{b}_0^{[l+1]}, \dots, \mathbf{W}_0^{[L]}, \mathbf{b}_0^{[L]} \right),
\end{aligned}$$

obviously $\tilde{\mathcal{T}}_{l,s}^\alpha$ is a linear operator, thus $\mathcal{T}_{l,s}^\alpha$ is an affine operator. \square

Directly from Lemma 4.1 and Lemma 4.2, we obtain that both one-step null embedding and one-step splitting embedding satisfy the property of output preserving and representation preserving, and all we need is to check the property of criticality preserving.

Proposition (Proposition 4.1 in main text). For any one-step null embedding $\mathcal{T}_{l,0}^\alpha$, given any NN $(\{m_l\}_{l=0}^L)$ and its parameters $\boldsymbol{\theta}_{\text{narr}} \in \text{Tuple}_{\{m_0, \dots, m_L\}}$ with $\text{Tuple}_{\{m_0, \dots, m_L\}} \in \mathcal{D}_{l,0}$, we have $\boldsymbol{\theta}_{\text{wide}} := \mathcal{T}_{l,0}^\alpha(\boldsymbol{\theta}_{\text{narr}})$ satisfies the following conditions: given any data S , loss $\ell(\cdot, \cdot)$ and activation $\sigma(\cdot)$, if $\nabla_{\boldsymbol{\theta}} R_S(\boldsymbol{\theta}_{\text{narr}}) = \mathbf{0}$, then $\nabla_{\boldsymbol{\theta}} R_S(\boldsymbol{\theta}_{\text{wide}}) = \mathbf{0}$.

Proof. Gradient of loss with respect to network parameters of each layer can be computed from F , G , and Z as follows

$$\begin{aligned}\nabla_{\mathbf{W}^{[l']}} R_S(\boldsymbol{\theta}) &= \nabla_{\mathbf{W}^{[l']}} \mathbb{E}_S \ell(\mathbf{f}_{\boldsymbol{\theta}}(\mathbf{x}), \mathbf{y}) = \mathbb{E}_S \left(\left(\mathbf{z}_{\boldsymbol{\theta}}^{[l']} \circ \mathbf{g}_{\boldsymbol{\theta}}^{[l']} \right) (\mathbf{f}_{\boldsymbol{\theta}}^{[l'-1]})^\top \right), \\ \nabla_{\mathbf{b}^{[l']}} R_S(\boldsymbol{\theta}) &= \nabla_{\mathbf{b}^{[l']}} \mathbb{E}_S \ell(\mathbf{f}_{\boldsymbol{\theta}}(\mathbf{x}), \mathbf{y}) = \mathbb{E}_S \left(\mathbf{z}_{\boldsymbol{\theta}}^{[l']} \circ \mathbf{g}_{\boldsymbol{\theta}}^{[l']} \right).\end{aligned}$$

Then, by Lemma 4.1, we have for $l' \neq l, l+1$,

$$\nabla_{\mathbf{W}^{[l']}} R_S(\boldsymbol{\theta}_{\text{wide}}) = \nabla_{\mathbf{W}^{[l']}} R_S(\boldsymbol{\theta}_{\text{narr}}) = \mathbf{0},$$

and

$$\nabla_{\mathbf{b}^{[l']}} R_S(\boldsymbol{\theta}_{\text{wide}}) = \nabla_{\mathbf{b}^{[l']}} R_S(\boldsymbol{\theta}_{\text{narr}}) = \mathbf{0}.$$

Also, for any $j \in [m_{l+1}]$, $k \in [m_l]$, since

$$(\mathbf{z}_{\boldsymbol{\theta}_{\text{wide}}}^{[l+1]})_j = (\mathbf{z}_{\boldsymbol{\theta}_{\text{narr}}}^{[l+1]})_j, \quad (\mathbf{g}_{\boldsymbol{\theta}_{\text{wide}}}^{[l+1]})_j = (\mathbf{g}_{\boldsymbol{\theta}_{\text{narr}}}^{[l+1]})_j, \quad (\mathbf{f}_{\boldsymbol{\theta}_{\text{wide}}}^{[l]})_k = (\mathbf{f}_{\boldsymbol{\theta}_{\text{narr}}}^{[l]})_k,$$

and $\mathbf{W}_{j, (m_{l+1})}^{[l+1]} \equiv 0$, we obtain that

$$\begin{aligned}\nabla_{\mathbf{W}_{j,k}^{[l+1]}} R_S(\boldsymbol{\theta}_{\text{wide}}) &= \nabla_{\mathbf{W}_{j,k}^{[l+1]}} R_S(\boldsymbol{\theta}_{\text{narr}}) = 0, \\ \nabla_{\mathbf{W}_{j, (m_{l+1})}^{[l+1]}} R_S(\boldsymbol{\theta}_{\text{wide}}) &= 0, \\ \nabla_{\mathbf{b}_j^{[l+1]}} R_S(\boldsymbol{\theta}_{\text{wide}}) &= \nabla_{\mathbf{b}_j^{[l+1]}} R_S(\boldsymbol{\theta}_{\text{narr}}) = 0.\end{aligned}$$

Similarly, for any $j \in [m_l]$, $k \in [m_{l-1}]$, since

$$(\mathbf{z}_{\boldsymbol{\theta}_{\text{wide}}}^{[l]})_j = (\mathbf{z}_{\boldsymbol{\theta}_{\text{narr}}}^{[l]})_j, \quad (\mathbf{g}_{\boldsymbol{\theta}_{\text{wide}}}^{[l]})_j = (\mathbf{g}_{\boldsymbol{\theta}_{\text{narr}}}^{[l]})_j, \quad (\mathbf{f}_{\boldsymbol{\theta}_{\text{wide}}}^{[l-1]})_k = (\mathbf{f}_{\boldsymbol{\theta}_{\text{narr}}}^{[l-1]})_k,$$

and $\mathbf{W}_{(m_{l+1}), k}^{[l]} \equiv 0$, we have

$$\begin{aligned}\nabla_{\mathbf{W}_{j,k}^{[l]}} R_S(\boldsymbol{\theta}_{\text{wide}}) &= \nabla_{\mathbf{W}_{j,k}^{[l]}} R_S(\boldsymbol{\theta}_{\text{narr}}) = 0, \\ \nabla_{\mathbf{W}_{(m_{l+1}), k}^{[l]}} R_S(\boldsymbol{\theta}_{\text{wide}}) &= 0, \\ \nabla_{\mathbf{b}_j^{[l]}} R_S(\boldsymbol{\theta}_{\text{wide}}) &= \nabla_{\mathbf{b}_j^{[l]}} R_S(\boldsymbol{\theta}_{\text{narr}}) = 0.\end{aligned}$$

Moreover, by Lemma 4.1, the output function $\mathbf{f}_{\boldsymbol{\theta}_{\text{wide}}}^{[L]} = \mathbf{f}_{\boldsymbol{\theta}_{\text{narr}}}^{[L]}$ is independent of the hyperparameter α , and $R_S(\boldsymbol{\theta}_{\text{wide}}) = R_S(\boldsymbol{\theta}_{\text{narr}})$, then since $\mathbf{b}_{(m_{l+1})}^{[l]} = \alpha$, we have

$$\nabla_{\mathbf{b}_{(m_{l+1})}^{[l]}} R_S(\boldsymbol{\theta}_{\text{wide}}) = \frac{\partial}{\partial \alpha} R_S(\boldsymbol{\theta}_{\text{narr}}) = 0.$$

Collecting all the above relations, we obtain that $\nabla_{\boldsymbol{\theta}} R_S(\boldsymbol{\theta}_{\text{wide}}) = \mathbf{0}$. □

Proposition (Proposition 4.2 in main text). For any one-step splitting embedding $\mathcal{T}_{l,s}^\alpha$, given any $\text{NN}(\{m_l\}_{l=0}^L)$ and its parameters $\boldsymbol{\theta}_{\text{narr}} \in \text{Tuple}_{\{m_0, \dots, m_L\}}$ with $\text{Tuple}_{\{m_0, \dots, m_L\}} \in \mathcal{D}_{l,s}$, we have $\boldsymbol{\theta}_{\text{wide}} := \mathcal{T}_{l,s}^\alpha(\boldsymbol{\theta}_{\text{narr}})$ satisfies the following conditions: given any data S , loss $\ell(\cdot, \cdot)$ and activation $\sigma(\cdot)$, if $\nabla_{\boldsymbol{\theta}} R_S(\boldsymbol{\theta}_{\text{narr}}) = \mathbf{0}$, then $\nabla_{\boldsymbol{\theta}} R_S(\boldsymbol{\theta}_{\text{wide}}) = \mathbf{0}$.

Proof. Gradient of loss with respect to network parameters of each layer can be computed from F , G , and Z as follows

$$\begin{aligned}\nabla_{\mathbf{W}^{[l']}} R_S(\boldsymbol{\theta}) &= \nabla_{\mathbf{W}^{[l']}} \mathbb{E}_S \ell(\mathbf{f}_{\boldsymbol{\theta}}(\mathbf{x}), \mathbf{y}) = \mathbb{E}_S \left(\left(\mathbf{z}_{\boldsymbol{\theta}}^{[l']} \circ \mathbf{g}_{\boldsymbol{\theta}}^{[l']} \right) (\mathbf{f}_{\boldsymbol{\theta}}^{[l'-1]})^\top \right), \\ \nabla_{\mathbf{b}^{[l']}} R_S(\boldsymbol{\theta}) &= \nabla_{\mathbf{b}^{[l']}} \mathbb{E}_S \ell(\mathbf{f}_{\boldsymbol{\theta}}(\mathbf{x}), \mathbf{y}) = \mathbb{E}_S \left(\mathbf{z}_{\boldsymbol{\theta}}^{[l']} \circ \mathbf{g}_{\boldsymbol{\theta}}^{[l']} \right).\end{aligned}$$

Then, by Lemma 4.2, we have for $l' \neq l, l+1$,

$$\nabla_{\mathbf{W}^{[l']}} R_S(\boldsymbol{\theta}_{\text{wide}}) = \nabla_{\mathbf{W}^{[l']}} R_S(\boldsymbol{\theta}_{\text{narr}}) = \mathbf{0},$$

and

$$\nabla_{\mathbf{b}^{[l']}} R_S(\boldsymbol{\theta}_{\text{wide}}) = \nabla_{\mathbf{b}^{[l']}} R_S(\boldsymbol{\theta}_{\text{narr}}) = \mathbf{0}.$$

Also, for any $j \in [m_{l+1}]$, $k \in [m_l]$, since

$$\begin{aligned}(\mathbf{z}_{\boldsymbol{\theta}_{\text{wide}}}^{[l+1]})_j &= (\mathbf{z}_{\boldsymbol{\theta}_{\text{narr}}}^{[l+1]})_j, & (\mathbf{g}_{\boldsymbol{\theta}_{\text{wide}}}^{[l+1]})_j &= (\mathbf{g}_{\boldsymbol{\theta}_{\text{narr}}}^{[l+1]})_j, \\ (\mathbf{f}_{\boldsymbol{\theta}_{\text{wide}}}^{[l]})_k &= (\mathbf{f}_{\boldsymbol{\theta}_{\text{narr}}}^{[l]})_k, & (\mathbf{f}_{\boldsymbol{\theta}_{\text{wide}}}^{[l]})_{m_l+1} &= (\mathbf{f}_{\boldsymbol{\theta}_{\text{narr}}}^{[l]})_s,\end{aligned}$$

we obtain

$$\begin{aligned}\nabla_{\mathbf{W}_{j,k}^{[l+1]}} R_S(\boldsymbol{\theta}_{\text{wide}}) &= \nabla_{\mathbf{W}_{j,k}^{[l+1]}} R_S(\boldsymbol{\theta}_{\text{narr}}) = 0, \\ \nabla_{\mathbf{W}_{j,(m_l+1)}^{[l+1]}} R_S(\boldsymbol{\theta}_{\text{wide}}) &= 0, \\ \nabla_{\mathbf{b}_j^{[l+1]}} R_S(\boldsymbol{\theta}_{\text{wide}}) &= \nabla_{\mathbf{b}_j^{[l+1]}} R_S(\boldsymbol{\theta}_{\text{narr}}) = 0.\end{aligned}$$

Similarly, for any $j \in [m_l] \setminus \{s\}$,

$$(\mathbf{z}_{\boldsymbol{\theta}_{\text{wide}}}^{[l]})_j = (\mathbf{z}_{\boldsymbol{\theta}_{\text{narr}}}^{[l]})_j, \quad (\mathbf{z}_{\boldsymbol{\theta}_{\text{wide}}}^{[l]})_s = (1 - \alpha)(\mathbf{z}_{\boldsymbol{\theta}_{\text{narr}}}^{[l]})_s, \quad (\mathbf{z}_{\boldsymbol{\theta}_{\text{wide}}}^{[l]})_{(m_l+1)} = \alpha(\mathbf{z}_{\boldsymbol{\theta}_{\text{narr}}}^{[l]})_s,$$

and for any $i \in [m_l]$,

$$(\mathbf{g}_{\boldsymbol{\theta}_{\text{wide}}}^{[l]})_i = (\mathbf{g}_{\boldsymbol{\theta}_{\text{narr}}}^{[l]})_i, \quad (\mathbf{g}_{\boldsymbol{\theta}_{\text{wide}}}^{[l]})_{(m_l+1)} = (\mathbf{g}_{\boldsymbol{\theta}_{\text{narr}}}^{[l]})_s,$$

and for $k \in [m_{l-1}]$,

$$(\mathbf{f}_{\boldsymbol{\theta}_{\text{wide}}}^{[l-1]})_k = (\mathbf{f}_{\boldsymbol{\theta}_{\text{narr}}}^{[l-1]})_k,$$

hence for any $j \in [m_l] \setminus \{s\}$, $k \in [m_{l-1}]$:

$$\begin{aligned}\nabla_{\mathbf{W}_{j,k}^{[l]}} R_S(\boldsymbol{\theta}_{\text{wide}}) &= \nabla_{\mathbf{W}_{j,k}^{[l]}} R_S(\boldsymbol{\theta}_{\text{narr}}) = 0, \\ \nabla_{\mathbf{b}_j^{[l]}} R_S(\boldsymbol{\theta}_{\text{wide}}) &= \nabla_{\mathbf{b}_j^{[l]}} R_S(\boldsymbol{\theta}_{\text{narr}}) = 0,\end{aligned}$$

$$\begin{aligned}
\nabla_{\mathbf{W}_{s,k}^{[l]}} R_S(\boldsymbol{\theta}_{\text{wide}}) &= (1 - \alpha) \nabla_{\mathbf{W}_{s,k}^{[l]}} R_S(\boldsymbol{\theta}_{\text{narr}}) = 0, \\
\nabla_{\mathbf{W}_{(m_j+1),k}^{[l]}} R_S(\boldsymbol{\theta}_{\text{wide}}) &= \alpha \nabla_{\mathbf{W}_{s,k}^{[l]}} R_S(\boldsymbol{\theta}_{\text{narr}}) = 0, \\
\nabla_{\mathbf{b}_s^{[l]}} R_S(\boldsymbol{\theta}_{\text{wide}}) &= (1 - \alpha) \nabla_{\mathbf{b}_s^{[l]}} R_S(\boldsymbol{\theta}_{\text{narr}}) = 0, \\
\nabla_{\mathbf{b}_{(m_l+1)}^{[l]}} R_S(\boldsymbol{\theta}_{\text{wide}}) &= \alpha \nabla_{\mathbf{b}_s^{[l]}} R_S(\boldsymbol{\theta}_{\text{narr}}) = 0.
\end{aligned}$$

Collecting all the above relations, we obtain that $\nabla_{\boldsymbol{\theta}} R_S(\boldsymbol{\theta}_{\text{wide}}) = \mathbf{0}$. \square

Combining altogether Lemma 4.1, Lemma 4.2, Proposition 4.1 and Proposition 4.2, we finish our proof for Theorem 4.1.

Theorem (Theorem 4.2 in main text). A K -step composition embedding is a critical embedding.

Proof. We shall prove it using induction.

For $K = 1$, Proposition 4.2 holds since both one-step null embedding and one step splitting embedding are critical embeddings.

Assume that Proposition 4.2 holds for $K = l - 1$, we want to show that it also holds for $K = l$.

From the induction hypothesis, we only need to show that if given two critical embeddings $\mathcal{T}_1, \mathcal{T}_2$, then $\mathcal{T}_2 \mathcal{T}_1$ is also a critical embedding.

(i) $\mathcal{T}_2 \mathcal{T}_1$ is injective:

For $\boldsymbol{\theta}_1$ and $\boldsymbol{\theta}_2$ belonging to a same tuple class but $\boldsymbol{\theta}_1 \neq \boldsymbol{\theta}_2$, since \mathcal{T}_1 is injective, then $\mathcal{T}_1(\boldsymbol{\theta}_1) \neq \mathcal{T}_1(\boldsymbol{\theta}_2)$. Since \mathcal{T}_2 is injective, then $\mathcal{T}_2 \mathcal{T}_1(\boldsymbol{\theta}_1) \neq \mathcal{T}_2 \mathcal{T}_1(\boldsymbol{\theta}_2)$.

(ii) $\mathcal{T}_2 \mathcal{T}_1$ is an affine embedding:

We use the fact that the composition of two affine operators is affine and hence we finish our proof.

(iii) $\mathcal{T}_2 \mathcal{T}_1$ satisfies the property of output preserving:

Since \mathcal{T}_1 satisfies the property of output preserving, then for any $\boldsymbol{\theta}$, $f_{\mathcal{T}_1}(\boldsymbol{\theta}) = f\boldsymbol{\theta}$. Similarly for \mathcal{T}_2 , we have $f_{\mathcal{T}_2 \mathcal{T}_1}(\boldsymbol{\theta}) = f_{\mathcal{T}_1}(\boldsymbol{\theta})$, hence $f_{\mathcal{T}_2 \mathcal{T}_1}(\boldsymbol{\theta}) = f\boldsymbol{\theta}$.

(iv) $\mathcal{T}_2 \mathcal{T}_1$ satisfies the property of representation preserving:

Similar reasoning in (iii).

(v) $\mathcal{T}_2 \mathcal{T}_1$ satisfies the property of criticality preserving:

Since $\boldsymbol{\theta}$ is a critical point of $R_S(\boldsymbol{\theta})$, so is $\mathcal{T}_1(\boldsymbol{\theta})$, and $\mathcal{T}_2 \mathcal{T}_1(\boldsymbol{\theta})$ as well, we finish the proof. \square

Lemma (Lemma 5.1 in main text). For any affine embedding $\mathcal{T} : \text{Tuple}_{\{m_0, \dots, m_L\}} \rightarrow \text{Tuple}_{\{m'_0, \dots, m'_L\}}$ satisfying the output preserving property, if there exists a total index mapping $\mathcal{I} = \{\mathcal{I}_l\}_{l=0}^L$ from $\text{NN}(\{m'_l\}_{l=0}^L)$ to $\text{NN}(\{m_l\}_{l=0}^L)$ and auxiliary variables $\boldsymbol{\beta} = \{\boldsymbol{\beta}_j^{[l]} \in \mathbb{R} \mid l \in [0 : L], j \in [m'_l] \setminus \mathcal{I}_l^{-1}(0)\}$, such that for any given neuron belonging to $\text{NN}(\{m'_l\}_{l=0}^L)$, located in layer l with index j , the following two statements hold:

(i) If $\mathcal{I}_l(j) \neq 0$, $(f_{\boldsymbol{\theta}_{\text{wide}}}^{[l]})_j = (f_{\boldsymbol{\theta}_{\text{narr}}}^{[l]})_{\mathcal{I}_l(j)}$ and $(e_{\boldsymbol{\theta}_{\text{wide}}}^{[l]})_j = \boldsymbol{\beta}_j^{[l]} (e_{\boldsymbol{\theta}_{\text{narr}}}^{[l]})_{\mathcal{I}_l(j)}$,

(ii) If $\mathcal{I}_l(j) = 0$, $(\mathbf{f}_{\boldsymbol{\theta}_{\text{wide}}}^{[l]})_j = \text{Const}$ and $(\mathbf{e}_{\boldsymbol{\theta}_{\text{wide}}}^{[l]})_j = 0$,

then \mathcal{T} is a critical embedding.

Proof. For any critical point $\boldsymbol{\theta}_{\text{narr}}^c \in \Theta_{\text{narr}}^c$, we set $\boldsymbol{\theta}_{\text{wide}} := \mathcal{T}(\boldsymbol{\theta}_{\text{narr}}^c)$. Then, since we have for any $l' \in [L]$

$$\begin{aligned}\nabla_{\mathbf{w}^{[l']}} R_S(\boldsymbol{\theta}) &= \nabla_{\mathbf{w}^{[l']}} \mathbb{E}_S \ell(\mathbf{f}_{\boldsymbol{\theta}}(x), \mathbf{y}) = \mathbb{E}_S \left(\left(\mathbf{z}_{\boldsymbol{\theta}}^{[l']} \circ \mathbf{g}_{\boldsymbol{\theta}}^{[l']} \right) (\mathbf{f}_{\boldsymbol{\theta}}^{[l'-1]})^\top \right), \\ \nabla_{\mathbf{b}^{[l']}} R_S(\boldsymbol{\theta}) &= \nabla_{\mathbf{b}^{[l']}} \mathbb{E}_S \ell(\mathbf{f}_{\boldsymbol{\theta}}(x), \mathbf{y}) = \mathbb{E}_S \left(\mathbf{z}_{\boldsymbol{\theta}}^{[l']} \circ \mathbf{g}_{\boldsymbol{\theta}}^{[l']} \right).\end{aligned}$$

Hence for $\mathcal{I}_l(i), \mathcal{I}_{l-1}(j) \neq 0$

$$\begin{aligned}\nabla_{(\mathbf{w}_{\text{wide}}^{[l]})_{ij}} R_S(\boldsymbol{\theta}_{\text{wide}}) &= \mathbb{E}_S (\mathbf{e}_{\boldsymbol{\theta}_{\text{wide}}}^{[l]})_i (\mathbf{f}_{\boldsymbol{\theta}_{\text{wide}}}^{[l-1]})_j \\ &= \mathbb{E}_S \boldsymbol{\beta}_i^{[l]} (\mathbf{e}_{\boldsymbol{\theta}_{\text{narr}}^c}^{[l]})_{\mathcal{I}_l(i)} (\mathbf{f}_{\boldsymbol{\theta}_{\text{narr}}^c}^{[l-1]})_{\mathcal{I}_{l-1}(j)} \\ &= \boldsymbol{\beta}_i^{[l]} \nabla_{(\mathbf{w}_{\text{narr}}^{[l]})_{\mathcal{I}_l(i), \mathcal{I}_{l-1}(j)}} R_S(\boldsymbol{\theta}_{\text{narr}}^c) = 0, \\ \nabla_{(\mathbf{b}_{\text{wide}}^{[l]})_i} R_S(\boldsymbol{\theta}_{\text{wide}}) &= \mathbb{E}_S (\mathbf{e}_{\boldsymbol{\theta}_{\text{wide}}}^{[l]})_i \\ &= \mathbb{E}_S \boldsymbol{\beta}_i^{[l]} (\mathbf{e}_{\boldsymbol{\theta}_{\text{narr}}^c}^{[l]})_{\mathcal{I}_l(i)} \\ &= \boldsymbol{\beta}_i^{[l]} \nabla_{(\mathbf{b}_{\text{narr}}^{[l]})_{\mathcal{I}_l(i)}} R_S(\boldsymbol{\theta}_{\text{narr}}^c) = 0.\end{aligned}$$

By condition (ii), these gradients are obviously 0 for $\mathcal{I}_l(i) = 0$ or $\mathcal{I}_{l-1}(j) = 0$. Therefore, $\boldsymbol{\theta}_{\text{wide}}$ is also a critical point and \mathcal{T} is criticality preserving.

Since \mathcal{I} is a total index mapping, for any feature vector of $\text{NN}(\{m_l\}_{l=0}^L)$, the component of which is also the output function of at least a neuron in the wide NN by condition (i). Moreover, any neuron output function of a neuron in the wide NN is either constant or output function of a neuron in the narrow NN by condition (i) and (ii). Therefore, \mathcal{T} is representation preserving. Then \mathcal{T} is a critical embedding. \square

Theorem (Theorem 5.1 in main text). General compatible embedding is a critical embedding.

Remark that we later name it as general compatible critical embedding in this work.

Proof. We need to prove four properties one by one.

1. $\mathcal{T}_{\mathcal{I}}^\alpha$ is output preserving and representation preserving;
2. $\mathcal{T}_{\mathcal{I}}^\alpha$ is an injective operator;
3. $\mathcal{T}_{\mathcal{I}}^\alpha$ is an affine embedding;
4. $\mathcal{T}_{\mathcal{I}}^\alpha$ is criticality preserving.

(i) We prove output preserving and representation preserving by doing induction on layers.

For the first layer, i.e., $l = 1$, we have for $i \in [m'_1]$,

$$\begin{aligned} (\mathbf{f}_{\boldsymbol{\theta}_{\text{wide}}}^{[1]})_i &= \sigma \left(\sum_{j \in [m'_0]} (\mathbf{W}_{\text{wide}}^{[1]})_{ij} \mathbf{x}_j + (\mathbf{b}_{\text{wide}}^{[1]})_i \right) \\ &= \sigma \left(\sum_{j \in [m_0]} \boldsymbol{\alpha}_{ij}^{[1]} (\mathbf{W}_{\text{narr}}^{[1]})_{\mathcal{I}_1(i),j} \mathbf{x}_j + (\boldsymbol{\alpha}_{\mathbf{b}}^{[1]})_i + (\mathbf{b}_{\text{narr}}^{[1]})_{\mathcal{I}_1(i)} \right), \end{aligned}$$

then for $i \notin \mathcal{I}_1^{-1}(0)$, for each $j \in [m_0]$, since $\boldsymbol{\alpha}_{ij}^{[1]} = \sum_{s \in \mathcal{I}_0^{-1}(j)} \boldsymbol{\alpha}_{is} = 1$, $(\boldsymbol{\alpha}_{\mathbf{b}}^{[1]})_i = 0$, hence

$$(\mathbf{f}_{\boldsymbol{\theta}_{\text{wide}}}^{[1]})_i = (\mathbf{f}_{\boldsymbol{\theta}_{\text{narr}}}^{[1]})_{\mathcal{I}_1(i)}.$$

Otherwise, for $i \in \mathcal{I}_1^{-1}(0)$, we have $\boldsymbol{\alpha}_{ij}^{[1]} = 0$ and $(\boldsymbol{\alpha}_{\mathbf{b}}^{[1]})_i = (\mathbf{b}_*^{[1]})_i$, then

$$(\mathbf{f}_{\boldsymbol{\theta}_{\text{wide}}}^{[1]})_i = \sigma \left((\mathbf{b}_*^{[1]})_i \right),$$

which is a constant function playing the same role as the bias term in the next layer.

Suppose for layer $l - 1$, we have for $i \notin \mathcal{I}_{l-1}^{-1}(0)$

$$(\mathbf{f}_{\boldsymbol{\theta}_{\text{wide}}}^{[l-1]})_i = (\mathbf{f}_{\boldsymbol{\theta}_{\text{narr}}}^{[l-1]})_{\mathcal{I}_{l-1}(i)},$$

and for $i \in \mathcal{I}_{l-1}^{-1}(0)$,

$$(\mathbf{f}_{\boldsymbol{\theta}_{\text{wide}}}^{[l-1]})_i = \sigma \left((\mathbf{b}_*^{[l-1]})_i \right).$$

Then we want to show that this is also the case for layer l .

We obtain that

$$\begin{aligned} (\mathbf{f}_{\boldsymbol{\theta}_{\text{wide}}}^{[l]})_i &= \sigma \left(\sum_{j \in [m'_{l-1}]} (\mathbf{W}_{\text{wide}}^{[l]})_{ij} (\mathbf{f}_{\boldsymbol{\theta}_{\text{wide}}}^{[l-1]})_j + (\mathbf{b}_{\text{wide}}^{[l]})_i \right) \\ &= \sigma \left(\sum_{s \in [m_{l-1}] \cup \{0\}} \sum_{j \in \mathcal{I}_{l-1}^{-1}(s)} \boldsymbol{\alpha}_{ij}^{[l]} (\mathbf{W}_{\text{narr}}^{[l]})_{\mathcal{I}_l(i),s} (\mathbf{f}_{\boldsymbol{\theta}_{\text{wide}}}^{[l-1]})_j + (\boldsymbol{\alpha}_{\mathbf{b}}^{[l]})_i + (\mathbf{b}_{\text{narr}}^{[l]})_{\mathcal{I}_l(i)} \right) \\ &= \sigma \left(\sum_{s \in [m_{l-1}]} (\mathbf{W}_{\text{narr}}^{[l]})_{\mathcal{I}_l(i),s} (\mathbf{f}_{\boldsymbol{\theta}_{\text{narr}}}^{[l-1]})_s \sum_{j \in \mathcal{I}_{l-1}^{-1}(s)} \boldsymbol{\alpha}_{ij}^{[l]} + \sum_{j \in \mathcal{I}_{l-1}^{-1}(0)} \boldsymbol{\alpha}_{ij}^{[l]} \sigma \left((\mathbf{b}_*^{[l-1]})_j \right) \right. \\ &\quad \left. + (\boldsymbol{\alpha}_{\mathbf{b}}^{[l]})_i + (\mathbf{b}_{\text{narr}}^{[l]})_{\mathcal{I}_l(i)} \right). \end{aligned}$$

For $i \notin \mathcal{I}_l^{-1}(0)$, we have

$$(\mathbf{f}_{\boldsymbol{\theta}_{\text{wide}}}^{[l]})_i = \sigma \left(\sum_{s \in [m_{l-1}]} (\mathbf{W}_{\text{narr}}^{[l]})_{\mathcal{I}_l(i),s} (\mathbf{f}_{\boldsymbol{\theta}_{\text{narr}}}^{[l-1]})_s + (\mathbf{b}_{\text{narr}}^{[l]})_{\mathcal{I}_l(i)} \right) = (\mathbf{f}_{\boldsymbol{\theta}_{\text{narr}}}^{[l]})_{\mathcal{I}_l(i)}.$$

Otherwise, for $i \in \mathcal{I}_l^{-1}(0)$, we have

$$(\mathbf{f}_{\boldsymbol{\theta}_{\text{wide}}}^{[l]})_i = \sigma \left((\mathbf{b}_*^{[l]})_i \right).$$

Then, for any layer l , we have for $i \notin \mathcal{I}_l^{-1}(0)$,

$$(\mathbf{f}_{\boldsymbol{\theta}_{\text{wide}}}^{[l]})_i = (\mathbf{f}_{\boldsymbol{\theta}_{\text{narr}}}^{[l]})_{\mathcal{I}_l(i)},$$

and for $i \in \mathcal{I}_l^{-1}(0)$,

$$(\mathbf{f}_{\boldsymbol{\theta}_{\text{wide}}}^{[l]})_i = \sigma \left((\mathbf{b}_*^{[l]})_i \right).$$

Hence, we have proved already that $\mathcal{T}_{\mathcal{I}}^{\alpha}$ is output preserving and representation preserving.

(ii) We prove that $\mathcal{T}_{\mathcal{I}}^{\alpha}$ is injective.

If for $\boldsymbol{\theta}_{\text{narr},1}, \boldsymbol{\theta}_{\text{narr},2} \in \text{Tuple}_{\{m_0, \dots, m_L\}}$ and $\boldsymbol{\theta}_{\text{narr},1} \neq \boldsymbol{\theta}_{\text{narr},2}$, then there exists some $l \in [L]$, such that $\mathbf{W}_{\text{narr},1}^{[l]} \neq \mathbf{W}_{\text{narr},2}^{[l]}$ or $\mathbf{b}_{\text{narr},1}^{[l]} \neq \mathbf{b}_{\text{narr},2}^{[l]}$.

Then, if $\mathcal{T}_{\mathcal{I}}^{\alpha}$ is not injective, there exists $\boldsymbol{\theta}_{\text{narr},1} \neq \boldsymbol{\theta}_{\text{narr},2}$, such that $\boldsymbol{\theta}_{\text{wide},1} = \boldsymbol{\theta}_{\text{wide},2}$, where $\boldsymbol{\theta}_{\text{wide},1} := \mathcal{T}_{\mathcal{I}}^{\alpha}(\boldsymbol{\theta}_{\text{narr},1})$ and $\boldsymbol{\theta}_{\text{wide},2} := \mathcal{T}_{\mathcal{I}}^{\alpha}(\boldsymbol{\theta}_{\text{narr},2})$, and we want to show that this will never happen.

Since there exists $l \in [L]$, such that $\mathbf{W}_{\text{narr},1}^{[l]} \neq \mathbf{W}_{\text{narr},2}^{[l]}$ or $\mathbf{b}_{\text{narr},1}^{[l]} \neq \mathbf{b}_{\text{narr},2}^{[l]}$. For the case $\mathbf{W}_{\text{narr},1}^{[l]} \neq \mathbf{W}_{\text{narr},2}^{[l]}$, we obtain that there exists $i \in [m_{l+1}], j \in [m_l]$, such that $(\mathbf{W}_{\text{narr},1}^{[l]})_{i,j} \neq (\mathbf{W}_{\text{narr},2}^{[l]})_{i,j}$. We observe that $\mathbf{W}_{\text{wide},1}^{[l]} = \boldsymbol{\alpha}^{[l]} \circ \mathbf{W}_{\text{inter},1}^{[l]}$, and $\mathbf{W}_{\text{wide},2}^{[l]} = \boldsymbol{\alpha}^{[l]} \circ \mathbf{W}_{\text{inter},2}^{[l]}$. Then, since \mathcal{I} is a total index mapping, then for $i, j \neq 0$, $\mathcal{I}_l^{-1}(i), \mathcal{I}_{l-1}^{-1}(j) \neq \emptyset$, hence for any $k \in \mathcal{I}_l^{-1}(i)$ and $l \in \mathcal{I}_{l-1}^{-1}(j)$,

$$\left(\mathbf{W}_{\text{wide},1}^{[l]} \right)_{k,l} = \boldsymbol{\alpha}_{k,l}^{[l]} \left(\mathbf{W}_{\text{narr},1}^{[l]} \right)_{i,j},$$

and

$$\left(\mathbf{W}_{\text{wide},2}^{[l]} \right)_{k,l} = \boldsymbol{\alpha}_{k,l}^{[l]} \left(\mathbf{W}_{\text{narr},2}^{[l]} \right)_{i,j}.$$

If $\mathbf{W}_{\text{wide},1}^{[l]} = \mathbf{W}_{\text{wide},2}^{[l]}$, then

$$\sum_{s \in \mathcal{I}_{l-1}^{-1}(j)} \left(\mathbf{W}_{\text{wide},1}^{[l]} \right)_{k,s} = \sum_{s \in \mathcal{I}_{l-1}^{-1}(j)} \left(\mathbf{W}_{\text{wide},2}^{[l]} \right)_{k,s},$$

hence

$$\sum_{s \in \mathcal{I}_l^{-1}(j)} \alpha_{k,s}^{[l]} \left(\mathbf{W}_{\text{narr},1}^{[l]} \right)_{i,j} = \sum_{s \in \mathcal{I}_l^{-1}(j)} \alpha_{k,s}^{[l]} \left(\mathbf{W}_{\text{narr},2}^{[l]} \right)_{i,j},$$

and since $\sum_{s \in \mathcal{I}_l^{-1}(j)} \alpha_{k,s}^{[l]} = 1$, we obtain that

$$\left(\mathbf{W}_{\text{narr},1}^{[l]} \right)_{i,j} = \left(\mathbf{W}_{\text{narr},2}^{[l]} \right)_{i,j},$$

which contradicts $\boldsymbol{\theta}_{\text{narr},1} \neq \boldsymbol{\theta}_{\text{narr},2}$.

For the case where $\mathbf{b}_{\text{narr},1}^{[l]} \neq \mathbf{b}_{\text{narr},2}^{[l]}$, since $(\alpha_{\mathbf{b}}^{[l]})_i = 0$ for any $l \in [L]$ with $i \notin \mathcal{I}_l^{-1}(0)$. Then for any $k \in \mathcal{I}_l^{-1}(j)$, $j \neq 0$,

$$\left(\mathbf{b}_{\text{wide},1}^{[l]} \right)_k = \alpha_{\mathbf{b}}^{[l]} + \left(\mathbf{b}_{\text{narr},1}^{[l]} \right)_j = \left(\mathbf{b}_{\text{narr},1}^{[l]} \right)_j,$$

and

$$\left(\mathbf{b}_{\text{wide},2}^{[l]} \right)_k = \alpha_{\mathbf{b}}^{[l]} + \left(\mathbf{b}_{\text{narr},2}^{[l]} \right)_j = \left(\mathbf{b}_{\text{narr},2}^{[l]} \right)_j,$$

hence $\mathbf{b}_{\text{wide},1}^{[l]} \neq \mathbf{b}_{\text{wide},2}^{[l]}$, and we finish the injection proof.

(iii) We prove that $\mathcal{T}_{\mathcal{I}}^{\alpha}$ is affine.

It is obvious that $\mathcal{T}_{\mathcal{I}}^{\alpha}$ is affine since for any $\boldsymbol{\theta} \in \text{Tuple}_{\{m_0, \dots, m_L\}}$, $\tilde{\mathcal{T}}_{\mathcal{I}}^{\alpha}(\boldsymbol{\theta}) := \mathcal{T}_{\mathcal{I}}^{\alpha}(\boldsymbol{\theta}) - \mathcal{T}_{\mathcal{I}}^{\alpha}(\mathbf{0})$ puts the weights and biases of null neuron as zero, and multiplies the weights and biases of effective by some constant, thus $\tilde{\mathcal{T}}_{\mathcal{I}}^{\alpha}$ is a linear operator.

(iv) We prove that $\mathcal{T}_{\mathcal{I}}^{\alpha}$ is criticality preserving.

We only need to check whether or not the conditions in Lemma 5.1 are satisfied. Moreover, we want to show that the collection of auxiliary variables

$$\boldsymbol{\beta} := \left\{ \beta_j^{[l]} \in \mathbb{R} \mid l \in [0 : L], j \in [m'_l] \setminus \mathcal{I}_l^{-1}(0) \right\}$$

in Condition 1 are exactly the auxiliary variables

$$\boldsymbol{\beta} = \left\{ \beta_j^{[l]} \in \mathbb{R} \mid l \in [0 : L], j \in [m'_l] \setminus \mathcal{I}_l^{-1}(0) \right\}$$

in Lemma 5.1. We show them using induction.

For layer L , by definition

$$\mathbf{e}_{\boldsymbol{\theta}_{\text{wide}}}^{[L]} = \mathbf{z}_{\boldsymbol{\theta}_{\text{wide}}}^{[L]} \circ \mathbf{g}_{\boldsymbol{\theta}_{\text{wide}}}^{[L]},$$

since $\mathcal{T}_{\mathcal{I}}^{\alpha}$ is output preserving for any activation, then

$$\mathbf{e}_{\boldsymbol{\theta}_{\text{wide}}}^{[L]} = \mathbf{z}_{\boldsymbol{\theta}_{\text{wide}}}^{[L]} \circ \mathbf{g}_{\boldsymbol{\theta}_{\text{wide}}}^{[L]} = \mathbf{z}_{\boldsymbol{\theta}_{\text{narr}}}^{[L]} \circ \mathbf{g}_{\boldsymbol{\theta}_{\text{narr}}}^{[L]} = \mathbf{e}_{\boldsymbol{\theta}_{\text{narr}}}^{[L]},$$

hence for $i \in [m_L]$, $(\mathbf{e}_{\boldsymbol{\theta}_{\text{wide}}}^{[L]})_i = \beta_i^{[L]} (\mathbf{e}_{\boldsymbol{\theta}_{\text{narr}}}^{[L]})_i$, since $\beta_i^{[L]} = 1$ for $i \in [m_L]$.

For any $l \in [L]$, suppose for layer l , we have for $i \notin \mathcal{I}_l^{-1}(0)$,

$$\left(\mathbf{e}_{\boldsymbol{\theta}_{\text{wide}}}^{[l]}\right)_i = \boldsymbol{\beta}_i^{[l]} \left(\mathbf{e}_{\boldsymbol{\theta}_{\text{narr}}}^{[l]}\right)_{\mathcal{I}_l(i)},$$

and for $i \in \mathcal{I}_l^{-1}(0)$,

$$\left(\mathbf{e}_{\boldsymbol{\theta}_{\text{wide}}}^{[l]}\right)_i = 0.$$

Then for layer $l-1$, we have

$$\begin{aligned} \left(\mathbf{e}_{\boldsymbol{\theta}_{\text{wide}}}^{[l-1]}\right)_j &= \left(\mathbf{g}_{\boldsymbol{\theta}_{\text{wide}}}^{[l-1]}\right)_j \sum_{i \in [m_l']} (\mathbf{W}_{\text{wide}}^{[l]})_{ij} \left(\mathbf{e}_{\boldsymbol{\theta}_{\text{wide}}}^{[l]}\right)_i \\ &= \left(\mathbf{g}_{\boldsymbol{\theta}_{\text{wide}}}^{[l-1]}\right)_j \sum_{k \in [m_l] \cup \{0\}} \sum_{i \in \mathcal{I}_l^{-1}(k)} \boldsymbol{\alpha}_{ij}^{[l]} (\mathbf{W}_{\text{inter}}^{[l]})_{ij} \left(\mathbf{e}_{\boldsymbol{\theta}_{\text{wide}}}^{[l]}\right)_i \\ &= \left(\mathbf{g}_{\boldsymbol{\theta}_{\text{wide}}}^{[l-1]}\right)_j \sum_{k \in [m_l]} (\mathbf{W}_{\text{narr}}^{[l]})_{k, \mathcal{I}_{l-1}(j)} \left(\mathbf{e}_{\boldsymbol{\theta}_{\text{narr}}}^{[l]}\right)_k \sum_{i \in \mathcal{I}_l^{-1}(k)} \boldsymbol{\alpha}_{ij}^{[l]} \boldsymbol{\beta}_i^{[l]}, \end{aligned}$$

from Condition 1, we observe that $\boldsymbol{\beta}_j^{[l-1]} := \sum_{i \in \mathcal{I}_l^{-1}(k)} \boldsymbol{\alpha}_{ij}^{[l]} \boldsymbol{\beta}_i^{[l]}$, then

$$\left(\mathbf{e}_{\boldsymbol{\theta}_{\text{wide}}}^{[l-1]}\right)_j = \left(\mathbf{g}_{\boldsymbol{\theta}_{\text{wide}}}^{[l-1]}\right)_j \boldsymbol{\beta}_j^{[l-1]} \sum_{k \in [m_l]} (\mathbf{W}_{\text{narr}}^{[l]})_{k, \mathcal{I}_{l-1}(j)} \left(\mathbf{e}_{\boldsymbol{\theta}_{\text{narr}}}^{[l]}\right)_k.$$

Since the choice of activation $\sigma(\cdot)$ is arbitrary, we may follow the same argument used for $\mathbf{f}_{\boldsymbol{\theta}_{\text{wide}}}^{[l]}$ and $\mathbf{f}_{\boldsymbol{\theta}_{\text{narr}}}^{[l]}$ in item (i), where for $i \notin \mathcal{I}_l^{-1}(0)$,

$$\left(\mathbf{f}_{\boldsymbol{\theta}_{\text{wide}}}^{[l]}\right)_i = \left(\mathbf{f}_{\boldsymbol{\theta}_{\text{narr}}}^{[l]}\right)_{\mathcal{I}_l(i)},$$

and for $i \in \mathcal{I}_l^{-1}(0)$,

$$\left(\mathbf{f}_{\boldsymbol{\theta}_{\text{wide}}}^{[l]}\right)_i = \sigma\left(\left(\mathbf{b}_*^{[l]}\right)_i\right).$$

Similarly, for $i \notin \mathcal{I}_l^{-1}(0)$,

$$\left(\mathbf{g}_{\boldsymbol{\theta}_{\text{wide}}}^{[l]}\right)_i = \left(\mathbf{g}_{\boldsymbol{\theta}_{\text{narr}}}^{[l]}\right)_{\mathcal{I}_l(i)},$$

and for $i \in \mathcal{I}_l^{-1}(0)$,

$$\left(\mathbf{g}_{\boldsymbol{\theta}_{\text{wide}}}^{[l]}\right)_i = \sigma^{(1)}\left(\left(\mathbf{b}_*^{[l]}\right)_i\right).$$

Hence for $j \notin \mathcal{I}_{l-1}^{-1}(0)$,

$$\begin{aligned} \left(\mathbf{e}_{\boldsymbol{\theta}_{\text{wide}}}^{[l-1]}\right)_j &= \left(\mathbf{g}_{\boldsymbol{\theta}_{\text{wide}}}^{[l-1]}\right)_j \boldsymbol{\beta}_j^{[l-1]} \sum_{k \in [m_l]} (\mathbf{W}_{\text{narr}}^{[l]})_{k, \mathcal{I}_{l-1}(j)} \left(\mathbf{e}_{\boldsymbol{\theta}_{\text{narr}}}^{[l]}\right)_k \\ &= \boldsymbol{\beta}_j^{[l-1]} \left(\mathbf{g}_{\boldsymbol{\theta}_{\text{narr}}}^{[l-1]}\right)_{\mathcal{I}_{l-1}(j)} \sum_{k \in [m_l]} (\mathbf{W}_{\text{narr}}^{[l]})_{k, \mathcal{I}_{l-1}(j)} \left(\mathbf{e}_{\boldsymbol{\theta}_{\text{narr}}}^{[l]}\right)_k \\ &= \boldsymbol{\beta}_j^{[l-1]} \left(\mathbf{e}_{\boldsymbol{\theta}_{\text{narr}}}^{[l-1]}\right)_j. \end{aligned}$$

Otherwise, for $j \in \mathcal{I}_{l-1}^{-1}(0)$, we set $\beta_j^{[l-1]} = 0$, then

$$(e_{\theta_{\text{wide}}}^{[l-1]})_j = \beta_j^{[l-1]} \sigma^{(1)} \left((\mathbf{b}_*^{[l-1]})_j \right) (z_{\theta_{\text{narr}}}^{[l-1]})_{\mathcal{I}_{l-1}(j)} = 0.$$

From the above proof, we find out that the variables

$$\boldsymbol{\beta} := \left\{ \beta_j^{[l]} \in \mathbb{R} \mid l \in [0 : L], j \in [m'_l] \setminus \mathcal{I}_l^{-1}(0) \right\}$$

in Condition 1 are exactly the variables in Lemma 5.1, hence we finish our proof. \square

Theorem (Theorem 6.2 in main text). Given an NN($\{m_l\}_{l=0}^L$) and any of its parameters $\boldsymbol{\theta} \in \mathbb{R}^M$, for any critical embedding $\mathcal{T} : \mathbb{R}^M \rightarrow \mathbb{R}^{M'}$ to any wider NN($\{m'_l\}_{l=0}^L$), the number of positive, zero, negative eigenvalues of $\mathbf{H}_S(\mathcal{T}(\boldsymbol{\theta}))$ is no less than the counterparts of $\mathbf{H}_S(\boldsymbol{\theta})$.

Proof. Because \mathcal{T} is a critical embedding, therefore, it is an affine injective operator associated with $\mathbf{A} \in \mathbb{R}^{M' \times M}, \mathbf{c} \in \mathbb{R}^{M'}$, such that $\mathcal{T}(\boldsymbol{\theta}) = \mathbf{A}\boldsymbol{\theta} + \mathbf{c}$. By the output preserving property of \mathcal{T} , we have

$$R_S(\boldsymbol{\theta}) \equiv R_S(\mathbf{A}\boldsymbol{\theta} + \mathbf{c}).$$

Hence,

$$\nabla_{\boldsymbol{\theta}} \nabla_{\boldsymbol{\theta}} R_S(\boldsymbol{\theta}) \equiv \nabla_{\boldsymbol{\theta}} \nabla_{\boldsymbol{\theta}} R_S(\mathbf{A}\boldsymbol{\theta} + \mathbf{c}).$$

Then

$$\mathbf{A}^\top \mathbf{H}_S(\mathbf{A}\boldsymbol{\theta} + \mathbf{c}) \mathbf{A} \equiv \mathbf{H}_S(\boldsymbol{\theta}).$$

Given any $\boldsymbol{\theta}_0$, if $\mathbf{H}_S(\boldsymbol{\theta}_0)$ has k negative eigenvalues $\{\lambda_j^{\text{neg}}\}_{j=1}^k$ with associated orthonormal eigenvectors $\{e_j^{\text{neg}}\}_{j=1}^k$, then $\{\mathbf{A}e_j^{\text{neg}}\}_{j=1}^k$ satisfies, for any e_j^{neg} ,

$$(\mathbf{A}e_j^{\text{neg}})^\top \mathbf{H}_S(\mathbf{A}\boldsymbol{\theta}_0 + \mathbf{c}) \mathbf{A}e_j^{\text{neg}} = (e_j^{\text{neg}})^\top \mathbf{H}_S(\boldsymbol{\theta}_0) e_j^{\text{neg}} = \lambda_j^{\text{neg}} < 0. \quad (\text{A.11})$$

By full rankness of \mathbf{A} , we have

$$\dim \left(\text{span} \left(\left\{ \mathbf{A}e_j^{\text{neg}} \right\}_{j=1}^k \right) \right) = k.$$

Thus, $\mathbf{H}_S(\mathbf{A}\boldsymbol{\theta}_0 + \mathbf{c})$ has at least k negative eigenvalues. Similarly, we can prove this result for the number of zero and positive eigenvalues.

Hence, in particular, for any critical embedding \mathcal{T} , the number of negative eigenvalues of $\mathbf{H}_S(\boldsymbol{\theta})$ is no more than the counterpart of $\mathbf{H}_S(\mathcal{T}(\boldsymbol{\theta}))$. \square

We would like to introduce some additional notations in order to state Lemma 6.1 and Lemma 6.2. In order to calculate the Hessian $\mathbf{H}_S(\boldsymbol{\theta}) = \nabla_{\boldsymbol{\theta}} \nabla_{\boldsymbol{\theta}} R_S(\boldsymbol{\theta})$, we need to compute $v_S(\boldsymbol{\theta})$:

$$v_S(\boldsymbol{\theta}) := \mathbb{E}_S \nabla \ell(\mathbf{f}(\mathbf{x}, \boldsymbol{\theta}), \mathbf{f}^*(\mathbf{x}))^\top \nabla_{\boldsymbol{\theta}} \mathbf{f}_{\boldsymbol{\theta}}(\mathbf{x}) = \sum_{i=1}^{m_L} \mathbb{E}_S \partial_i \ell(\mathbf{f}_{\boldsymbol{\theta}}, \mathbf{f}^*) \nabla_{\boldsymbol{\theta}} (\mathbf{f}_{\boldsymbol{\theta}})_i,$$

where $\partial_i \ell(\mathbf{f}_\theta, \mathbf{f}^*)$ is the i -th element of $\nabla \ell(\mathbf{f}(\mathbf{x}, \theta), \mathbf{f}^*(\mathbf{x}))$, and $(\mathbf{f}_\theta)_i$ is the i -th element of vector \mathbf{f}_θ , then for the Hessian $\mathbf{H}_S(\theta)$, we have

$$\begin{aligned} \mathbf{H}_S(\theta) &= \nabla_\theta \nabla_\theta R_S(\theta) \\ &= \sum_{i=1}^{m_L} \mathbb{E}_S \nabla_\theta (\partial_i \ell(\mathbf{f}_\theta, \mathbf{f}^*)) \nabla_\theta (\mathbf{f}_\theta)_i + \sum_{i=1}^{m_L} \mathbb{E}_S \partial_i \ell(\mathbf{f}_\theta, \mathbf{f}^*) \nabla_\theta \nabla_\theta ((\mathbf{f}_\theta)_i) \\ &= \sum_{i,j=1}^{m_L} \mathbb{E}_S \partial_{ij} \ell(\mathbf{f}_\theta, \mathbf{f}^*) \nabla_\theta (\mathbf{f}_\theta)_i (\nabla_\theta (\mathbf{f}_\theta)_j)^\top + \sum_{i=1}^{m_L} \mathbb{E}_S \partial_i \ell(\mathbf{f}_\theta, \mathbf{f}^*) \nabla_\theta \nabla_\theta ((\mathbf{f}_\theta)_i), \end{aligned}$$

where $\partial_{ij} \ell(\mathbf{f}_\theta, \mathbf{f}^*)$ is the (i, j) -th element of $\nabla \nabla \ell(\mathbf{f}(\mathbf{x}, \theta), \mathbf{f}^*(\mathbf{x}))$.

Hence $\mathbf{H}_S^{(1)}(\theta)$ and $\mathbf{H}_S^{(2)}(\theta)$ respectively becomes

$$\mathbf{H}_S^{(1)}(\theta) := \sum_{i,j=1}^{m_L} \mathbb{E}_S \partial_{ij} \ell(\mathbf{f}_\theta, \mathbf{f}^*) \nabla_\theta (\mathbf{f}_\theta)_i (\nabla_\theta (\mathbf{f}_\theta)_j)^\top, \quad (\text{A.12})$$

and

$$\mathbf{H}_S^{(2)}(\theta) := \sum_{i=1}^{m_L} \mathbb{E}_S \partial_i \ell(\mathbf{f}_\theta, \mathbf{f}^*) \nabla_\theta \nabla_\theta ((\mathbf{f}_\theta)_i). \quad (\text{A.13})$$

Obviously, we observe that

$$\mathbf{H}_S(\theta) = \mathbf{H}_S^{(1)}(\theta) + \mathbf{H}_S^{(2)}(\theta). \quad (\text{A.14})$$

We observe that for $i \in [m_L], j \in [m_{L-1}]$,

$$(\mathbf{f}_\theta)_i = \sum_{j=1}^{m_{L-1}} (\mathbf{W}^{[L]})_{ij} (\mathbf{f}_\theta^{[L-1]})_j + (\mathbf{b}^{[L]})_i,$$

hence we obtain that, for any $i \in [m_L]$,

$$\begin{aligned} \frac{\partial (\mathbf{f}_\theta)_i}{\partial \mathbf{W}^{[L]}} &= \begin{bmatrix} \mathbf{0}^{(i-1) \times m_{L-1}} \\ (\mathbf{f}_\theta^{[L-1]})_j \\ \mathbf{0}^{(m_L-i) \times m_{L-1}} \end{bmatrix}, \quad j \in [m_{L-1}], \\ \frac{\partial (\mathbf{f}_\theta)_i}{\partial \mathbf{b}^{[L]}} &= \begin{bmatrix} \mathbf{0}^{(i-1) \times 1} \\ 1 \\ \mathbf{0}^{(m_L-i) \times 1} \end{bmatrix}. \end{aligned} \quad (\text{A.15})$$

Finally, given an NN($\{m_l\}_{l=0}^L$) and its parameter

$$\theta = (\mathbf{W}^{[1]}, \mathbf{b}^{[1]}, \dots, \mathbf{W}^{[L]}, \mathbf{b}^{[L]}) \in \text{Tuple}_{\{m_0, \dots, m_L\}},$$

we remind the readers once again that the collection of parameters θ is a $2L$ -tuple. We denote that the upper bracket $[L-1]$ by limiting ourselves to the first $2L-2$ element of the tuple, i.e.,

$$\theta^{[L-1]} := (\mathbf{W}^{[1]}, \mathbf{b}^{[1]}, \dots, \mathbf{W}^{[L-2]}, \mathbf{b}^{[L-2]}, \mathbf{W}^{[L-1]}, \mathbf{b}^{[L-1]}) \in \text{Tuple}_{\{m_0, \dots, m_{L-2}, m_{L-1}\}}, \quad (\text{A.16})$$

and similarly, we identify $\boldsymbol{\theta}^{[L-1]}$ with its vectorization $\text{vec}(\boldsymbol{\theta}^{[L-1]}) \in \mathbb{R}^{M^{[L-1]}}$ with $M^{[L-1]} = \sum_{l=0}^{L-2} (m_l + 1)m_{l+1}$. Then, we denote hereafter the expressions

$$\begin{aligned} \mathbf{H}_S^{(2),[L-1]}(\boldsymbol{\theta}) &:= \sum_{i,j=1}^{m_L} \mathbb{E}_S \partial_{ij} \ell(\mathbf{f}_{\boldsymbol{\theta}}, \mathbf{f}^*) \mathbf{W}_{ij}^{[L]} \nabla_{\boldsymbol{\theta}^{[L-1]}} \nabla_{\boldsymbol{\theta}^{[L-1]}} \left(\mathbf{f}_{\boldsymbol{\theta}}^{[L-1]} \right)_j, \\ \mathbf{H}_S^{(1),[L-1]}(\boldsymbol{\theta}) &:= \sum_{i,j=1}^{m_L} \mathbb{E}_S \partial_{ij} \ell(\mathbf{f}_{\boldsymbol{\theta}}, \mathbf{f}^*) \nabla_{\boldsymbol{\theta}^{[L-1]}}(\mathbf{f}_{\boldsymbol{\theta}})_i \left(\nabla_{\boldsymbol{\theta}^{[L-1]}}(\mathbf{f}_{\boldsymbol{\theta}})_j \right)^\top, \end{aligned}$$

and $\mathbf{H}_S^{[L-1]}(\boldsymbol{\theta}) := \mathbf{H}_S^{(1),[L-1]}(\boldsymbol{\theta}) + \mathbf{H}_S^{(2),[L-1]}(\boldsymbol{\theta})$.

Lemma (Lemma 6.1 in main text). Given any data S , loss $\ell(\cdot, \cdot)$ and activation $\sigma(\cdot)$, for any NN, if a critical point $\boldsymbol{\theta}^c \in \boldsymbol{\Theta}^c$ satisfies:

- (i) $\mathbf{H}_S(\boldsymbol{\theta}^c) \succeq 0$;
- (ii) $\mathbf{H}_S^{(2),[L-1]}(\boldsymbol{\theta}^c) \neq \mathbf{0}$,

then there exists a general compatible critical embedding \mathcal{T} , such that $\mathcal{T}(\boldsymbol{\theta}^c)$ is a strict-saddle point.

Proof. We consider the three-fold global splitting embedding $\mathcal{T}_{\text{global}}$ defined in Example 5.1. Obviously, $\mathcal{T}_{\text{global}}$ is a general compatible critical embedding. After choosing a critical point $\boldsymbol{\theta}_{\text{narr}}^c \in \boldsymbol{\Theta}^c$, and we have that $\boldsymbol{\theta}_{\text{wide}}^c := \mathcal{T}_{\text{global}}(\boldsymbol{\theta}_{\text{narr}}^c)$.

$\boldsymbol{\theta}_{\text{narr}}^c$ and $\boldsymbol{\theta}_{\text{wide}}^c$ are tuples, and we misuse these notations and identify them with their vectorizations, i.e., $\boldsymbol{\theta}_{\text{narr}}^c \in \mathbb{R}^M$ and $\boldsymbol{\theta}_{\text{wide}}^c \in \mathbb{R}^{M'}$ for some M and M' . More specifically, we set

$$\begin{aligned} \boldsymbol{\theta}_{\text{narr}}^c &= \left(\text{vec}(\mathbf{W}^{[1]})^\top, \dots, \text{vec}(\mathbf{W}^{[L-2]})^\top, \text{vec}(\mathbf{W}^{[L-1]})^\top, \right. \\ &\quad \left. \mathbf{b}^{[1]\top}, \dots, \mathbf{b}^{[L-1]\top}, \text{vec}(\mathbf{W}^{[L]})^\top, \mathbf{b}^{[L]\top} \right)^\top, \\ \boldsymbol{\theta}_{\text{wide}}^c &= \left(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \boldsymbol{\theta}_3, \mathbf{0}_{1 \times (M' - 3M + 2m_L)}, \text{vec}(\mathbf{W}^{[L]})^\top, \text{vec}(\mathbf{W}^{[L]})^\top, \text{vec}(-\mathbf{W}^{[L]})^\top, \mathbf{b}^{[L]\top} \right)^\top, \end{aligned}$$

with

$$\begin{aligned} \boldsymbol{\theta}_1 &:= \boldsymbol{\theta}_2 := \boldsymbol{\theta}_3 := \boldsymbol{\theta}^{[L-1]} \\ &= \left(\text{vec}(\mathbf{W}^{[1]})^\top, \dots, \text{vec}(\mathbf{W}^{[L-2]})^\top, \text{vec}(\mathbf{W}^{[L-1]})^\top, \mathbf{b}^{[1]\top}, \dots, \mathbf{b}^{[L-1]\top} \right)^\top, \end{aligned}$$

then we observe that

$$\boldsymbol{\theta}_{\text{narr}}^c = \left(\boldsymbol{\theta}_1, \text{vec}(\mathbf{W}^{[L]})^\top, \mathbf{b}^{[L]\top} \right)^\top,$$

and we are able to do some computations:

$$\begin{aligned} \mathbf{H}_S^{[L-1]}(\boldsymbol{\theta}_{\text{wide}}^c) &= \sum_{i,j=1}^{m_L} \mathbb{E}_S \partial_{ij} \ell(\mathbf{f}_{\boldsymbol{\theta}_{\text{wide}}}, \mathbf{f}^*) \nabla_{\boldsymbol{\theta}_{\text{wide}}^{[L-1]}}(\mathbf{f}_{\boldsymbol{\theta}_{\text{wide}}})_i \left(\nabla_{\boldsymbol{\theta}_{\text{wide}}^{[L-1]}}(\mathbf{f}_{\boldsymbol{\theta}_{\text{wide}}})_j \right)^\top \\ &\quad + \sum_{i,j=1}^{m_L} \mathbb{E}_S \partial_{ij} \ell(\mathbf{f}_{\boldsymbol{\theta}_{\text{wide}}}, \mathbf{f}^*) \left(\mathbf{W}_{\text{wide}}^{[L]} \right)_{ij} \nabla_{\boldsymbol{\theta}_{\text{wide}}^{[L-1]}} \nabla_{\boldsymbol{\theta}_{\text{wide}}^{[L-1]}} \left(\mathbf{f}_{\boldsymbol{\theta}_{\text{wide}}}^{[L-1]} \right)_j, \quad (\text{A.17}) \end{aligned}$$

Then for the first part

$$\sum_{i,j=1}^{m_L} \mathbb{E}_S \partial_{ij} \ell(\mathbf{f}_{\boldsymbol{\theta}_{\text{wide}}}, \mathbf{f}^*) \nabla_{\boldsymbol{\theta}_{\text{wide}}^{[L-1]}} (\mathbf{f}_{\boldsymbol{\theta}_{\text{wide}}})_i \left(\nabla_{\boldsymbol{\theta}_{\text{wide}}^{[L-1]}} (\mathbf{f}_{\boldsymbol{\theta}_{\text{wide}}})_j \right)^\top,$$

we have

$$\sum_{i,j=1}^{m_L} \mathbb{E}_S \partial_{ij} \ell(\mathbf{f}_{\boldsymbol{\theta}_{\text{wide}}}, \mathbf{f}^*) \nabla_{\boldsymbol{\theta}_{\text{wide}}^{[L-1]}} (\mathbf{f}_{\boldsymbol{\theta}_{\text{wide}}})_i \left(\nabla_{\boldsymbol{\theta}_{\text{wide}}^{[L-1]}} (\mathbf{f}_{\boldsymbol{\theta}_{\text{wide}}})_j \right)^\top = \begin{pmatrix} \mathbf{A}_{1,1} & \mathbf{A}_{1,2} & \mathbf{A}_{1,3} \\ \mathbf{A}_{2,1} & \mathbf{A}_{2,2} & \mathbf{A}_{2,3} \\ \mathbf{A}_{3,1} & \mathbf{A}_{3,2} & \mathbf{A}_{3,3} \end{pmatrix},$$

with

$$\mathbf{A}_{p,q} := \mathbb{E}_S \sum_{i,j=1}^{m_L} \partial_{ij} \ell(\mathbf{f}_{\boldsymbol{\theta}_{\text{wide}}}, \mathbf{f}^*) \nabla_{\boldsymbol{\theta}_p} (\mathbf{f}_{\boldsymbol{\theta}_{\text{wide}}})_i \left(\nabla_{\boldsymbol{\theta}_q} (\mathbf{f}_{\boldsymbol{\theta}_{\text{wide}}})_j \right)^\top,$$

for $p \in [3]$ and $q \in [3]$.

For the second part

$$\sum_{i,j=1}^{m_L} \mathbb{E}_S \partial_i \ell(\mathbf{f}_{\boldsymbol{\theta}_{\text{wide}}}, \mathbf{f}^*) \mathbf{W}_{ij}^{[L]} \nabla_{\boldsymbol{\theta}_{\text{wide}}^{[L-1]}} \nabla_{\boldsymbol{\theta}_{\text{wide}}^{[L-1]}} (\mathbf{f}_{\boldsymbol{\theta}_{\text{wide}}})_j,$$

we have

$$\sum_{i,j=1}^{m_L} \mathbb{E}_S \partial_i \ell(\mathbf{f}_{\boldsymbol{\theta}_{\text{wide}}}, \mathbf{f}^*) \mathbf{W}_{ij}^{[L]} \nabla_{\boldsymbol{\theta}_{\text{wide}}^{[L-1]}} \nabla_{\boldsymbol{\theta}_{\text{wide}}^{[L-1]}} (\mathbf{f}_{\boldsymbol{\theta}_{\text{wide}}})_j = \begin{pmatrix} \mathbf{B}_{1,1} & \mathbf{B}_{1,2} & \mathbf{B}_{1,3} \\ \mathbf{B}_{2,1} & \mathbf{B}_{2,2} & \mathbf{B}_{2,3} \\ \mathbf{B}_{3,1} & \mathbf{B}_{3,2} & \mathbf{B}_{3,3} \end{pmatrix},$$

with

$$\mathbf{B}_{p,q} := \mathbb{E}_S \sum_{i,j=1}^{m_L} \partial_i \ell(\mathbf{f}_{\boldsymbol{\theta}_{\text{wide}}}, \mathbf{f}^*) \mathbf{W}_{ij}^{[L]} \nabla_{\boldsymbol{\theta}_p} \nabla_{\boldsymbol{\theta}_q} (\mathbf{f}_{\boldsymbol{\theta}_{\text{wide}}})_j,$$

for $p \in [3]$ and $q \in [3]$.

Moreover,

$$\nabla_{\boldsymbol{\theta}_1} (\mathbf{f}_{\boldsymbol{\theta}_{\text{wide}}}(x))_i = \sum_{j=1}^{m_L-1} (\mathbf{W}_{\text{wide}}^{[L]})_{ij} \nabla_{\boldsymbol{\theta}_1} (\mathbf{f}_{\boldsymbol{\theta}_{\text{wide}}})_j, \quad i \in [m_L], \quad (\text{A.18a})$$

$$\nabla_{\boldsymbol{\theta}_2} (\mathbf{f}_{\boldsymbol{\theta}_{\text{wide}}}(x))_i = \sum_{j=m_L-1+1}^{2m_L-1} (\mathbf{W}_{\text{wide}}^{[L]})_{ij} \nabla_{\boldsymbol{\theta}_2} (\mathbf{f}_{\boldsymbol{\theta}_{\text{wide}}})_j, \quad i \in [m_L], \quad (\text{A.18b})$$

$$\nabla_{\boldsymbol{\theta}_3} (\mathbf{f}_{\boldsymbol{\theta}_{\text{wide}}}(x))_i = \sum_{j=2m_L-1+1}^{3m_L-1} (\mathbf{W}_{\text{wide}}^{[L]})_{ij} \nabla_{\boldsymbol{\theta}_3} (\mathbf{f}_{\boldsymbol{\theta}_{\text{wide}}})_j, \quad i \in [m_L], \quad (\text{A.18c})$$

hence by construction of $\mathcal{T}_{\text{global}}$, we obtain that

$$\nabla_{\boldsymbol{\theta}_1} (\mathbf{f}_{\boldsymbol{\theta}_{\text{wide}}}(x))_i = \nabla_{\boldsymbol{\theta}_2} (\mathbf{f}_{\boldsymbol{\theta}_{\text{wide}}}(x))_i = -\nabla_{\boldsymbol{\theta}_3} (\mathbf{f}_{\boldsymbol{\theta}_{\text{wide}}}(x))_i.$$

And for the Hessian, we have

$$\nabla_{\theta_1} \nabla_{\theta_1} (f_{\theta_{\text{wide}}}(\mathbf{x}))_i = \sum_{j=1}^{m_{L-1}} (\mathbf{W}_{\text{wide}}^{[L]})_{i,j} \nabla_{\theta_1} \nabla_{\theta_1} (f_{\theta_{\text{wide}}}^{[L-1]})_j, \quad i \in [m_L], \quad (\text{A.19a})$$

$$\nabla_{\theta_2} \nabla_{\theta_2} (f_{\theta_{\text{wide}}}(\mathbf{x}))_i = \sum_{j=m_{L-1}+1}^{2m_{L-1}} (\mathbf{W}_{\text{wide}}^{[L]})_{i,j} \nabla_{\theta_2} \nabla_{\theta_2} (f_{\theta_{\text{wide}}}^{[L-1]})_j, \quad i \in [m_L], \quad (\text{A.19b})$$

$$\nabla_{\theta_3} \nabla_{\theta_3} (f_{\theta_{\text{wide}}}(\mathbf{x}))_i = \sum_{j=2m_{L-1}+1}^{3m_{L-1}} (\mathbf{W}_{\text{wide}}^{[L]})_{i,j} \nabla_{\theta_3} \nabla_{\theta_3} (f_{\theta_{\text{wide}}}^{[L-1]})_j, \quad i \in [m_L], \quad (\text{A.19c})$$

$$\begin{aligned} \nabla_{\theta_1} \nabla_{\theta_2} (f_{\theta_{\text{wide}}}(\mathbf{x}))_i &= \nabla_{\theta_1} \nabla_{\theta_3} (f_{\theta_{\text{wide}}}(\mathbf{x}))_i \\ &= \nabla_{\theta_2} \nabla_{\theta_3} (f_{\theta_{\text{wide}}}(\mathbf{x}))_i = \mathbf{0}, \quad i \in [m_L], \end{aligned} \quad (\text{A.19d})$$

Thus for $\mathbf{H}_S^{(2),[L-1]}(\theta_{\text{narr}}^c) \neq \mathbf{0}$, then there exists a nonzero eigenvalue $\lambda \neq 0$ associated with its unit eigenvector \mathbf{v} .

For the cases where $\lambda > 0$, then $\mathbf{v}^\top \mathbf{H}_S^{(2),[L-1]}(\theta_{\text{narr}}^c) \mathbf{v} = \lambda$. We observe that the matrix $\tilde{\mathbf{H}}(\theta_{\text{wide}}^c)$ below is a principle submatrix of the Hessian at θ_{wide}^c , i.e., by choosing the columns and rows corresponding to $(\theta_1, \theta_2, \theta_3)$, we obtain that

$$\begin{aligned} \tilde{\mathbf{H}}(\theta_{\text{wide}}^c) := & \begin{bmatrix} \mathbf{H}_S^{(1),[L-1]}(\theta_{\text{narr}}^c) & \mathbf{H}_S^{(1),[L-1]}(\theta_{\text{narr}}^c) & -\mathbf{H}_S^{(1),[L-1]}(\theta_{\text{narr}}^c) \\ \mathbf{H}_S^{(1),[L-1]}(\theta_{\text{narr}}^c) & \mathbf{H}_S^{(1),[L-1]}(\theta_{\text{narr}}^c) & -\mathbf{H}_S^{(1),[L-1]}(\theta_{\text{narr}}^c) \\ -\mathbf{H}_S^{(1),[L-1]}(\theta_{\text{narr}}^c) & -\mathbf{H}_S^{(1),[L-1]}(\theta_{\text{narr}}^c) & \mathbf{H}_S^{(1),[L-1]}(\theta_{\text{narr}}^c) \end{bmatrix} \\ & + \begin{bmatrix} \mathbf{H}_S^{(2),[L-1]}(\theta_{\text{narr}}^c) & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{H}_S^{(2),[L-1]}(\theta_{\text{narr}}^c) & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & -\mathbf{H}_S^{(2),[L-1]}(\theta_{\text{narr}}^c) \end{bmatrix}, \end{aligned}$$

Then for $\mathbf{u} = [\frac{1}{2}\mathbf{v}^\top, \frac{1}{2}\mathbf{v}^\top, \mathbf{v}^\top]^\top$,

$$\mathbf{u}^\top \tilde{\mathbf{H}}(\theta_{\text{wide}}^c) \mathbf{u} = -\frac{1}{2}\lambda < 0,$$

indicating θ_{wide}^c is a strict-saddle point.

Otherwise, if $\lambda < 0$, since $\mathbf{v}^\top \tilde{\mathbf{H}}(\theta_{\text{wide}}^c) \mathbf{v} = \lambda$, then for $\mathbf{u} = [\mathbf{v}^\top, -\mathbf{v}^\top, \mathbf{0}]^\top$,

$$\mathbf{u}^\top \tilde{\mathbf{H}}(\theta_{\text{wide}}^c) \mathbf{u} = 2\lambda < 0,$$

indicating θ_{wide}^c is also a strict-saddle point. \square

Lemma (Lemma 6.2 in main text). Given any data S , loss $\ell(\cdot, \cdot)$ and activation $\sigma(\cdot)$, for any NN, if a critical point $\theta^c \in \Theta^c$ satisfies:

(i) $\mathbf{H}_S(\theta^c) \succeq \mathbf{0}$;

(ii) $\mathbf{H}_S^{(2),[L-1]}(\theta^c) = \mathbf{0}$;

(iii) $\mathbf{H}_S^{(2)}(\boldsymbol{\theta}^c) \neq \mathbf{0}$.

Then there exists a general compatible critical embedding \mathcal{T} , such that $\mathcal{T}(\boldsymbol{\theta}^c)$ is a strict-saddle point.

Proof. Since we have

$$\mathbf{H}_S^{(2)}(\boldsymbol{\theta}) := \sum_{i=1}^{m_L} \mathbb{E}_S \partial_i \ell(\mathbf{f}_{\boldsymbol{\theta}}, \mathbf{f}^*) \nabla_{\boldsymbol{\theta}} \nabla_{\boldsymbol{\theta}} ((\mathbf{f}_{\boldsymbol{\theta}})_i), \quad (\text{A.20})$$

then $\mathbf{H}_S^{(2)}(\boldsymbol{\theta}_{\text{narr}}^c)$ can be written into the form of

$$\mathbf{H}_S^{(2)}(\boldsymbol{\theta}_{\text{narr}}^c) = \mathbb{E}_S \begin{pmatrix} \mathbf{C}_{1,1} & \mathbf{C}_{1,2} & \mathbf{C}_{1,3} \\ \mathbf{C}_{2,1} & \mathbf{C}_{2,2} & \mathbf{C}_{2,3} \\ \mathbf{C}_{3,1} & \mathbf{C}_{3,2} & \mathbf{C}_{3,3} \end{pmatrix},$$

where

$$\begin{aligned} \mathbf{C}_{1,1} &= \sum_{i=1}^{m_L} \partial_i \ell(\mathbf{f}_{\boldsymbol{\theta}_{\text{narr}}}, \mathbf{f}^*) \nabla_{\mathbf{W}^{[L]}} \nabla_{\mathbf{W}^{[L]}} ((\mathbf{f}_{\boldsymbol{\theta}_{\text{narr}}})_i), \\ \mathbf{C}_{1,2} &= \sum_{i=1}^{m_L} \partial_i \ell(\mathbf{f}_{\boldsymbol{\theta}_{\text{narr}}}, \mathbf{f}^*) \nabla_{\mathbf{W}^{[L]}} \nabla_{\mathbf{b}^{[L]}} ((\mathbf{f}_{\boldsymbol{\theta}_{\text{narr}}})_i), \\ \mathbf{C}_{1,3} &= \sum_{i=1}^{m_L} \partial_i \ell(\mathbf{f}_{\boldsymbol{\theta}_{\text{narr}}}, \mathbf{f}^*) \nabla_{\mathbf{W}^{[L]}} \nabla_{\boldsymbol{\theta}^{[L-1]}} ((\mathbf{f}_{\boldsymbol{\theta}_{\text{narr}}})_i), \\ \mathbf{C}_{2,1} &= \sum_{i=1}^{m_L} \partial_i \ell(\mathbf{f}_{\boldsymbol{\theta}_{\text{narr}}}, \mathbf{f}^*) \nabla_{\mathbf{b}^{[L]}} \nabla_{\mathbf{W}^{[L]}} ((\mathbf{f}_{\boldsymbol{\theta}_{\text{narr}}})_i), \\ \mathbf{C}_{2,2} &= \sum_{i=1}^{m_L} \partial_i \ell(\mathbf{f}_{\boldsymbol{\theta}_{\text{narr}}}, \mathbf{f}^*) \nabla_{\mathbf{b}^{[L]}} \nabla_{\mathbf{b}^{[L]}} ((\mathbf{f}_{\boldsymbol{\theta}_{\text{narr}}})_i), \\ \mathbf{C}_{2,3} &= \sum_{i=1}^{m_L} \partial_i \ell(\mathbf{f}_{\boldsymbol{\theta}_{\text{narr}}}, \mathbf{f}^*) \nabla_{\mathbf{b}^{[L]}} \nabla_{\boldsymbol{\theta}^{[L-1]}} ((\mathbf{f}_{\boldsymbol{\theta}_{\text{narr}}})_i), \\ \mathbf{C}_{3,1} &= \sum_{i=1}^{m_L} \partial_i \ell(\mathbf{f}_{\boldsymbol{\theta}_{\text{narr}}}, \mathbf{f}^*) \nabla_{\boldsymbol{\theta}^{[L-1]}} \nabla_{\mathbf{W}^{[L]}} ((\mathbf{f}_{\boldsymbol{\theta}_{\text{narr}}})_i), \\ \mathbf{C}_{3,2} &= \sum_{i=1}^{m_L} \partial_i \ell(\mathbf{f}_{\boldsymbol{\theta}_{\text{narr}}}, \mathbf{f}^*) \nabla_{\boldsymbol{\theta}^{[L-1]}} \nabla_{\mathbf{b}^{[L]}} ((\mathbf{f}_{\boldsymbol{\theta}_{\text{narr}}})_i), \\ \mathbf{C}_{3,3} &= \sum_{i=1}^{m_L} \partial_i \ell(\mathbf{f}_{\boldsymbol{\theta}_{\text{narr}}}, \mathbf{f}^*) \nabla_{\boldsymbol{\theta}^{[L-1]}} \nabla_{\boldsymbol{\theta}^{[L-1]}} ((\mathbf{f}_{\boldsymbol{\theta}_{\text{narr}}})_i), \end{aligned}$$

from relation (A.15), we obtain that

$$\begin{aligned} & \mathbf{H}_S^{(2)}(\boldsymbol{\theta}_{\text{narr}}^c) \\ = & \mathbb{E}_S \begin{bmatrix} \mathbf{0} & \mathbf{0} & \sum_{i=1}^{m_L} \partial_i \ell(\mathbf{f}_{\boldsymbol{\theta}_{\text{narr}}}, \mathbf{f}^*) \nabla_{\mathbf{W}^{[L]}} \nabla_{\boldsymbol{\theta}^{[L-1]}} ((\mathbf{f}_{\boldsymbol{\theta}_{\text{narr}}})_i) \\ \sum_{i=1}^{m_L} \partial_i \ell(\mathbf{f}_{\boldsymbol{\theta}_{\text{narr}}}, \mathbf{f}^*) \nabla_{\boldsymbol{\theta}^{[L-1]}} \nabla_{\mathbf{W}^{[L]}} ((\mathbf{f}_{\boldsymbol{\theta}_{\text{narr}}})_i) & \mathbf{0} & \mathbf{0} \\ \sum_{i=1}^{m_L} \partial_i \ell(\mathbf{f}_{\boldsymbol{\theta}_{\text{narr}}}, \mathbf{f}^*) \nabla_{\boldsymbol{\theta}^{[L-1]}} \nabla_{\mathbf{b}^{[L]}} ((\mathbf{f}_{\boldsymbol{\theta}_{\text{narr}}})_i) & \mathbf{0} & \mathbf{H}_S^{(2),[L-1]}(\boldsymbol{\theta}_{\text{narr}}^c) \end{bmatrix}, \end{aligned}$$

We observe that the matrix $\tilde{\mathbf{H}}_S^{(2)}(\boldsymbol{\theta}_{\text{narr}}^c)$ below is a principle submatrix of the Hessian at $\boldsymbol{\theta}_{\text{narr}}^c$, i.e., by choosing the columns and rows corresponding to $(\mathbf{W}^{[L]}, \boldsymbol{\theta}^{[L-1]})$, we obtain that

$$\begin{aligned} & \tilde{\mathbf{H}}_S^{(2)}(\boldsymbol{\theta}_{\text{narr}}^c) \\ = & \mathbb{E}_S \begin{bmatrix} \mathbf{0} & \sum_{i=1}^{m_L} \partial_i \ell(\mathbf{f}_{\boldsymbol{\theta}_{\text{wide}}}, \mathbf{f}^*) \nabla_{\mathbf{W}^{[L]}} \nabla_{\boldsymbol{\theta}^{[L-1]}} ((\mathbf{f}_{\boldsymbol{\theta}_{\text{narr}}})_i) \\ \sum_{i=1}^{m_L} \partial_i \ell(\mathbf{f}_{\boldsymbol{\theta}_{\text{wide}}}, \mathbf{f}^*) \nabla_{\boldsymbol{\theta}^{[L-1]}} \nabla_{\mathbf{W}^{[L]}} ((\mathbf{f}_{\boldsymbol{\theta}_{\text{narr}}})_i) & \mathbf{H}_S^{(2), [L-1]}(\boldsymbol{\theta}_{\text{narr}}^c) \end{bmatrix}, \end{aligned}$$

and it satisfies $\text{tr}(\tilde{\mathbf{H}}_S^{(2)}(\boldsymbol{\theta}_{\text{narr}}^c)) = 0$.

Since $\mathbf{H}_S^{(2)}(\boldsymbol{\theta}_{\text{narr}}^c) \neq \mathbf{0}$, then $\tilde{\mathbf{H}}_S^{(2)}(\boldsymbol{\theta}_{\text{narr}}^c)$ has at least one negative eigenvalue, denoted by λ^{neg} associated with its unit eigenvector \mathbf{v} , i.e., $\mathbf{v}^\top \tilde{\mathbf{H}}_S^{(2)}(\boldsymbol{\theta}_{\text{narr}}^c) \mathbf{v} = \lambda^{\text{neg}} < 0$. Then we consider the three-fold global splitting embedding $\mathcal{T}_{\text{global}}$ defined in Example 5.1, with $\boldsymbol{\theta}_{\text{wide}}^c = \mathcal{T}_{\text{global}}(\boldsymbol{\theta}_{\text{narr}}^c)$.

From relation (A.14), we obtain that

$$\mathbf{H}_S(\boldsymbol{\theta}_{\text{wide}}^c) = \mathbf{H}_S^{(1)}(\boldsymbol{\theta}_{\text{wide}}^c) + \mathbf{H}_S^{(2)}(\boldsymbol{\theta}_{\text{wide}}^c),$$

we recall that

$$\boldsymbol{\theta}_1 = \boldsymbol{\theta}_2 = \boldsymbol{\theta}^{[L-1]} = \left(\text{vec}(\mathbf{W}^{[1]})^\top, \dots, \text{vec}(\mathbf{W}^{[L-2]})^\top, \text{vec}(\mathbf{W}^{[L-1]})^\top, \mathbf{b}^{[1]\top}, \dots, \mathbf{b}^{[L-1]\top} \right),$$

and if we concatenate $\boldsymbol{\theta}_1$ and $\text{vec}(\mathbf{W}^{[L]})^\top$ into a new vector $\mathbf{d}_1 := (\boldsymbol{\theta}_1, \text{vec}(\mathbf{W}^{[L]})^\top)$, and $\boldsymbol{\theta}_2$ and $\text{vec}(\mathbf{W}^{[L]})^\top$ into a new vector $\mathbf{d}_2 := (\boldsymbol{\theta}_2, \text{vec}(\mathbf{W}^{[L]})^\top)$, we observe that the matrix $\tilde{\mathbf{H}}(\boldsymbol{\theta}_{\text{wide}}^c)$ below is a principle submatrix of the Hessian $\mathbf{H}_S(\boldsymbol{\theta})$ at $\boldsymbol{\theta}_{\text{wide}}^c$, i.e., by choosing the columns and rows corresponding to $(\mathbf{d}_1, \mathbf{d}_2)$, we obtain that

$$\tilde{\mathbf{H}}(\boldsymbol{\theta}_{\text{wide}}^c) := \begin{bmatrix} \mathbf{D}_{1,1} & \mathbf{D}_{1,2} \\ \mathbf{D}_{2,1} & \mathbf{D}_{2,2} \end{bmatrix} + \begin{bmatrix} \tilde{\mathbf{H}}_S^{(2)}(\boldsymbol{\theta}_{\text{narr}}^c) & \mathbf{0} \\ \mathbf{0} & \tilde{\mathbf{H}}_S^{(2)}(\boldsymbol{\theta}_{\text{narr}}^c) \end{bmatrix},$$

with

$$\mathbf{D}_{p,q} := \mathbb{E}_S \sum_{i,j=1}^{m_L} \partial_{ij} \ell(\mathbf{f}_{\boldsymbol{\theta}_{\text{wide}}}, \mathbf{f}^*) \nabla_{d_p} (\mathbf{f}_{\boldsymbol{\theta}_{\text{wide}}})_i \left(\nabla_{d_q} (\mathbf{f}_{\boldsymbol{\theta}_{\text{wide}}})_j \right)^\top,$$

for $p \in [2]$ and $q \in [2]$. Since $\mathbf{d}_1 = \mathbf{d}_2$, we obtain that

$$\mathbf{D}_{1,1} = \mathbf{D}_{2,2} = \mathbf{D}_{1,2} = \mathbf{D}_{2,1}.$$

Then, by choosing $\mathbf{u} = [\mathbf{v}^\top, -\mathbf{v}^\top]^\top$, we have

$$\mathbf{u}^\top \tilde{\mathbf{H}}(\boldsymbol{\theta}_{\text{wide}}^c) \mathbf{u} = 2\lambda^{\text{neg}} < 0,$$

indicating $\boldsymbol{\theta}_{\text{wide}}^c$ is a strict-saddle point. □

B Step by step backpropagation

In this section, we derive all the relations in (3.5) concerning $\{\mathbf{f}_\theta^{[l]}\}_{l=1}^L$, $\{\mathbf{g}_\theta^{[l]}\}_{l=1}^L$, $\{\mathbf{z}_\theta^{[l]}\}_{l=1}^L$. We shall write out all the relations once again

$$\mathbf{z}_\theta^{[L]} = \nabla \ell, \quad (\text{B.1})$$

$$\mathbf{z}_\theta^{[l]} = (\mathbf{W}^{[l+1]})^\top \left(\mathbf{z}_\theta^{[l+1]} \circ \mathbf{g}_\theta^{[l+1]} \right), \quad l \in [L-1], \quad (\text{B.2})$$

$$\nabla_{\mathbf{W}^{[l]}} \ell = \left(\mathbf{z}_\theta^{[l]} \circ \mathbf{g}_\theta^{[l]} \right) (\mathbf{f}_\theta^{[l-1]})^\top, \quad l \in [L], \quad (\text{B.3})$$

$$\nabla_{\mathbf{b}^{[l]}} \ell = \mathbf{z}_\theta^{[l]} \circ \mathbf{g}_\theta^{[l]}, \quad l \in [L]. \quad (\text{B.4})$$

For relation (B.1), since by definition $\mathbf{z}_\theta^{[L]} = \nabla_{\mathbf{f}^{[L]}} \ell = \nabla_{\mathbf{f}_\theta} \ell$, and $\ell = \ell(\mathbf{f}_\theta, \mathbf{f}^*)$, where \mathbf{f}^* is fixed, hence

$$\left(\mathbf{z}_\theta^{[L]} \right)_i = \frac{\partial \ell(\mathbf{f}_\theta, \mathbf{f}^*)}{\partial (\mathbf{f}_\theta)_i} = \partial_i \ell(\mathbf{f}_\theta, \mathbf{f}^*).$$

We prove relation (B.2) in detail and the rest of the relations follow the same reasoning. First things first, we show that it holds for $l = L-1$. Since $\mathbf{z}_\theta^{[L]} = \nabla \ell$, and $\mathbf{g}_\theta^{[L]} := \mathbf{1}$, hence $\mathbf{z}_\theta^{[L]} \circ \mathbf{g}_\theta^{[L]} = \nabla \ell$. By definition,

$$\mathbf{z}_\theta^{[L-1]} = \nabla_{\mathbf{f}^{[L-1]}} \ell,$$

using chain rule, since

$$(\mathbf{f}^{[L]})_i = \sum_{j=1}^{m_{L-1}} \mathbf{W}_{ij}^{[L]} (\mathbf{f}^{[L-1]})_j + (\mathbf{b}^{[L]})_i,$$

then

$$\left(\mathbf{z}_\theta^{[L-1]} \right)_j = \frac{\partial \ell}{\partial (\mathbf{f}^{[L-1]})_j} = \sum_{i=1}^{m_L} \frac{\partial \ell}{\partial (\mathbf{f}^{[L]})_i} \frac{\partial (\mathbf{f}^{[L]})_i}{\partial (\mathbf{f}^{[L-1]})_j} = \sum_{i=1}^{m_L} \partial_i \ell \mathbf{W}_{ij}^{[L]},$$

hence

$$\mathbf{z}_\theta^{[L-1]} = (\mathbf{W}^{[L]})^\top \nabla \ell = (\mathbf{W}^{[L]})^\top \left(\mathbf{z}_\theta^{[L]} \circ \mathbf{g}_\theta^{[L]} \right).$$

Now for $l \in [L-2]$, we have that since

$$\left(\mathbf{f}^{[l+1]} \right)_i = \sigma \left(\left(\sum_{j=1}^{m_l} \mathbf{W}_{ij}^{[l+1]} (\mathbf{f}^{[l]})_j \right) + (\mathbf{b}^{[l+1]})_i \right),$$

then

$$\left(\mathbf{z}_\theta^{[l]} \right)_j = \frac{\partial \ell}{\partial (\mathbf{f}^{[l]})_j} = \sum_{i=1}^{m_{l+1}} \frac{\partial \ell}{\partial (\mathbf{f}^{[l+1]})_i} \frac{\partial (\mathbf{f}^{[l+1]})_i}{\partial (\mathbf{f}^{[l]})_j},$$

where

$$\frac{\partial \left(\mathbf{f}^{[l+1]} \right)_i}{\partial \left(\mathbf{f}^{[l]} \right)_j} = \sigma^{(1)} \left(\left(\sum_{j=1}^{m_l} \mathbf{W}_{i,j}^{[l+1]} \left(\mathbf{f}^{[l]} \right)_j \right) + \left(\mathbf{b}^{[l+1]} \right)_i \right) \mathbf{W}_{i,j}^{[l+1]},$$

hence

$$\begin{aligned} \left(\mathbf{z}_\theta^{[l]} \right)_j &= \frac{\partial \ell}{\partial \mathbf{f}_j^{[l]}} = \sum_{i=1}^{m_{l+1}} \frac{\partial \ell}{\partial \mathbf{f}_i^{[l+1]}} \sigma^{(1)} \left(\left(\sum_{j=1}^{m_l} \mathbf{W}_{i,j}^{[l+1]} \mathbf{f}_j^{[l]} \right) + \mathbf{b}_i^{[l+1]} \right) \mathbf{W}_{i,j}^{[l+1]} \\ &= \left(\mathbf{W}^{[l+1]} \right)^\top \left(\mathbf{z}_\theta^{[l+1]} \circ \mathbf{g}_\theta^{[l+1]} \right). \end{aligned}$$

C Formal definition of the composition of two embeddings, Definition 4.5

Definition (Composition of two embeddings). Suppose we have an NN($\{m_l\}_{l=0}^L$) and its parameters $\theta \in \text{Tuple}_{\{m_0, \dots, m_L\}}$, and we have two embeddings \mathcal{T} and \mathcal{T}' satisfying

$$\mathcal{T} : \text{Tuple}_{\{m_0, \dots, m_L\}} \rightarrow \text{Tuple}_{\{m'_0, \dots, m'_L\}}, \quad \mathcal{T}' : \text{Tuple}_{\{m'_0, \dots, m'_L\}} \rightarrow \text{Tuple}_{\{m''_0, \dots, m''_L\}},$$

with \mathcal{T}' maps the range of \mathcal{T} into $\text{Tuple}_{\{m''_0, \dots, m''_L\}}$, where $m''_0 = m'_0 = m_0$, $m''_L = m'_L = m_L$, and for any $l \in [L-1]$, $m''_l \geq m'_l \geq m_l$. Since $\mathcal{T}'\mathcal{T}$ is obviously an injective operator, then $\mathcal{T}'\mathcal{T}$ is an embedding $\mathcal{T}'\mathcal{T} : \text{Tuple}_{\{m_0, \dots, m_L\}} \rightarrow \text{Tuple}_{\{m''_0, \dots, m''_L\}}$, and we term $\mathcal{T}'\mathcal{T}$ the composition of \mathcal{T}' and \mathcal{T} , i.e., for any $\theta \in \text{Tuple}_{\{m_0, \dots, m_L\}}$

$$\mathcal{T}'\mathcal{T}(\theta) := \mathcal{T}'(\mathcal{T}(\theta)).$$

References

- [1] Weinan E, Chao Ma, Lei Wu, and Stephan Wojtowytsch. Towards a mathematical understanding of neural network-based machine learning: What we know and what we don't. *CSIAM Transactions on Applied Mathematics*, 1(4):561–615, 2020.
- [2] Ruoyu Sun, Dawei Li, Shiyu Liang, Tian Ding, and Rayadurgam Srikant. The global landscape of neural networks: An overview. *IEEE Signal Processing Magazine*, 37(5):95–108, 2020.
- [3] Hao Li, Zheng Xu, Gavin Taylor, Christoph Studer, and Tom Goldstein. Visualizing the loss landscape of neural nets. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, pages 6391–6401, 2018.
- [4] Ivan Skorokhodov and Mikhail Burtsev. Loss landscape sightseeing with multi-point optimization. *arXiv preprint arXiv:1910.03867*, 2019.
- [5] Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and generalization in neural networks. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, pages 8580–8589, 2018.
- [6] Simon Du, Jason Lee, Haochuan Li, Liwei Wang, and Xiyu Zhai. Gradient descent finds global minima of deep neural networks. In *International Conference on Machine Learning*, pages 1675–1685. PMLR, 2019.

- [7] Sanjeev Arora, Simon S Du, Wei Hu, Zhiyuan Li, Ruslan Salakhutdinov, and Ruosong Wang. On exact computation with an infinitely wide neural net. In *Proceedings of the 33rd International Conference on Neural Information Processing Systems*, pages 8141–8150, 2019.
- [8] Lenaïc Chizat, Edouard Oyallon, and Francis Bach. On lazy training in differentiable programming. *Advances in Neural Information Processing Systems*, 32, 2019.
- [9] Difan Zou, Yuan Cao, Dongruo Zhou, and Quanguan Gu. Stochastic gradient descent optimizes over-parameterized deep relu networks. *arXiv preprint arXiv:1811.08888*, 2018.
- [10] Lénaïc Chizat and Francis R. Bach. On the global convergence of gradient descent for over-parameterized models using optimal transport. In *Advances in Neural Information Processing Systems*, 2018.
- [11] Tao Luo, Zhi-Qin John Xu, Zheng Ma, and Yaoyu Zhang. Phase diagram for two-layer relu neural networks at infinite-width limit. *Journal of Machine Learning Research*, 22(71):1–47, 2021.
- [12] Leo Breiman. Reflections after refereeing papers for nips. *The Mathematics of Generalization*, XX:11–15, 1995.
- [13] Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning requires rethinking generalization. In *5th International Conference on Learning Representations, ICLR 2017*.
- [14] Zhi-Qin J Xu, Yaoyu Zhang, and Yanyang Xiao. Training behavior of deep neural network in frequency domain. *International Conference on Neural Information Processing*, pages 264–274, 2019.
- [15] Zhi-Qin John Xu, Yaoyu Zhang, Tao Luo, Yanyang Xiao, and Zheng Ma. Frequency principle: Fourier analysis sheds light on deep neural networks. *Communications in Computational Physics*, 28(5):1746–1767, 2020.
- [16] Nasim Rahaman, Devansh Arpit, Aristide Baratin, Felix Draxler, Min Lin, Fred A Hamprecht, Yoshua Bengio, and Aaron Courville. On the spectral bias of deep neural networks. *International Conference on Machine Learning*, 2019.
- [17] Yaoyu Zhang, Tao Luo, Zheng Ma, and Zhi-Qin John Xu. A linear frequency principle model to understand the absence of overfitting in neural networks. *Chinese Physics Letters*, 38(3):038701, 2021.
- [18] Tao Luo, Zheng Ma, Zhi-Qin John Xu, and Yaoyu Zhang. Theory of the frequency principle for general deep neural networks. *CSIAM Transactions on Applied Mathematics*, 2(3):484–507, 2021.
- [19] Zhi-Qin John Xu, Yaoyu Zhang, and Tao Luo. Overview frequency principle/spectral bias in deep learning. *arXiv preprint arXiv:2201.07395*, 2022.
- [20] Hartmut Maennel, Olivier Bousquet, and Sylvain Gelly. Gradient descent quantizes relu network features. *arXiv preprint arXiv:1803.08367*, 2018.
- [21] Jonathan Frankle and Michael Carbin. The lottery ticket hypothesis: Finding sparse, trainable neural networks. In *International Conference on Learning Representations*, 2018.
- [22] Jason D Lee, Ioannis Panageas, Georgios Piliouras, Max Simchowitz, Michael I Jordan, and Benjamin Recht. First-order methods almost always avoid strict saddle points. *Mathematical Programming*, 176(1):311–337, 2019.
- [23] Yaim Cooper. Global minima of overparameterized neural networks. *SIAM Journal on Mathematics of Data Science*, 3(2):676–691, 2021.
- [24] Yaoyu Zhang, Zhongwang Zhang, Tao Luo, and Zhi-Qin John Xu. Embedding principle of loss landscape of deep neural networks. *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, 2021.
- [25] Kenji Fukumizu and Shun-ichi Amari. Local minima and plateaus in hierarchical structures of multilayer perceptrons. *Neural Networks*, 13(3):317–327, 2000.
- [26] Kenji Fukumizu, Shoichiro Yamaguchi, Yoh-ichi Mototake, and Mirai Tanaka. Semi-flat minima and saddle points by embedding neural networks to overparameterization. *Advances in Neural Information Processing Systems*, 32:13868–13876, 2019.

- [27] Berfin Simsek, François Ged, Arthur Jacot, Francesco Spadaro, Clement Hongler, Wulfram Gerstner, and Johanni Brea. Geometry of the loss landscape in overparameterized neural networks: Symmetries and invariances. In *Proceedings of the 38th International Conference on Machine Learning*, pages 9722–9732. PMLR, 2021.
- [28] Simon Du and Jason Lee. On the power of over-parametrization in neural networks with quadratic activation. In *International Conference on Machine Learning*, pages 1329–1338. PMLR, 2018.
- [29] Mahdi Soltanolkotabi, Adel Javanmard, and Jason D Lee. Theoretical insights into the optimization landscape of over-parameterized shallow neural networks. *IEEE Transactions on Information Theory*, 65(2):742–769, 2018.
- [30] Patrick Cheridito, Arnulf Jentzen, and Florian Rossmannek. Landscape analysis for shallow relu neural networks: Complete classification of critical points for affine target functions. *arXiv preprint arXiv:2103.10922*, 2021.
- [31] Nitish Shirish Keskar, Jorge Nocedal, Ping Tak Peter Tang, Dheevatsa Mudigere, and Mikhail Smelyanskiy. On large-batch training for deep learning: Generalization gap and sharp minima. In *5th International Conference on Learning Representations, ICLR 2017*, 2017.
- [32] Zhongwang Zhang, Hanxu Zhou, and Zhi-Qin John Xu. A variance principle explains why dropout finds flatter minima. *arXiv preprint arXiv:2111.01022*, 2021.
- [33] Lei Wu, Zhanxing Zhu, and Weinan E. Towards understanding generalization of deep learning: Perspective of loss landscapes. *arXiv preprint arXiv:1706.10239*, 2017.
- [34] Haowei He, Gao Huang, and Yang Yuan. Asymmetric valleys: beyond sharp and flat local minima. In *Proceedings of the 33rd International Conference on Neural Information Processing Systems*, pages 2553–2564, 2019.
- [35] Tian Ding, Dawei Li, and Ruoyu Sun. Sub-optimal local minima exist for neural networks with almost all non-linear activations. *arXiv preprint arXiv:1911.01413*, 2019.
- [36] Luca Venturi, Afonso S Bandeira, and Joan Bruna. Spurious valleys in one-hidden-layer neural network optimization landscapes. *Journal of Machine Learning Research*, 20:133, 2019.
- [37] Devansh Arpit, Stanisław Jastrzebski, Nicolas Ballas, David Krueger, Emmanuel Bengio, Maxinder S Kanwal, Tegan Maharaj, Asja Fischer, Aaron Courville, Yoshua Bengio, et al. A closer look at memorization in deep networks. In *International Conference on Machine Learning*, pages 233–242. PMLR, 2017.
- [38] Dimitris Kalimeris, Gal Kaplun, Preetum Nakkiran, Benjamin Edelman, Tristan Yang, Boaz Barak, and Haofeng Zhang. Sgd on neural networks learns functions of increasing complexity. *Advances in Neural Information Processing Systems*, 32:3496–3506, 2019.
- [39] Sebastian Goldt, Marc Mézard, Florent Krzakala, and Lenka Zdeborová. Modeling the influence of data structure on learning in neural networks: The hidden manifold model. *Physical Review X*, 10(4):041044, 2020.
- [40] Shilin He, Xing Wang, Shuming Shi, Michael R Lyu, and Zhaopeng Tu. Assessing the bilingual knowledge learned by neural machine translation models. *arXiv preprint arXiv:2004.13270*, 2020.
- [41] Chris Mingard, Joar Skalse, Guillermo Valle-Pérez, David Martínez-Rubio, Vladimir Mikulik, and Ard A Louis. Neural networks are a priori biased towards boolean functions with low entropy. *arXiv preprint arXiv:1909.11522*, 2019.
- [42] Pengzhan Jin, Lu Lu, Yifa Tang, and George Em Karniadakis. Quantifying the generalization error in deep learning in terms of data distribution and neural network smoothness. *Neural Networks*, 130:85–99, 2020.
- [43] Franco Pellegrini and Giulio Biroli. An analytic theory of shallow networks dynamics for hinge loss classification. *Advances in Neural Information Processing Systems*, 33, 2020.
- [44] Zhi-Qin John Xu, Hanxu Zhou, Tao Luo, and Yaoyu Zhang. Towards understanding the condensation of two-layer neural networks at initial training. *arXiv preprint arXiv:2105.11686*, 2021.

- [45] Anna Choromanska, Mikael Henaff, Michael Mathieu, Gérard Ben Arous, and Yann LeCun. The loss surfaces of multilayer networks. In *Artificial Intelligence and Statistics*, pages 192–204. PMLR, 2015.
- [46] Levent Sagun, Léon Bottou, and Yann LeCun. Singularity of the Hessian in deep learning. *arXiv preprint arXiv:1611.07476*, 2016.
- [47] Roger A Horn and Charles R Johnson. *Matrix Analysis*. Cambridge University Press, 2012.