

---

# Rank Diminishing in Deep Neural Networks

---

**Ruili Feng<sup>1</sup>, Kecheng Zheng<sup>2,1</sup>, Yukun Huang<sup>1</sup>, Deli Zhao<sup>2,3</sup>,**  
**Michael Jordan<sup>4</sup>, Zheng-Jun Zha<sup>1</sup>**

<sup>1</sup>University of Science and Technology of China, Hefei, China

<sup>2</sup>Ant Research, <sup>3</sup>Alibaba Group, Hangzhou, China

<sup>4</sup>University of California, Berkeley

[ruilifengustc@gmail.com](mailto:ruilifengustc@gmail.com), [{zkcys001,kevinh}@mail.ustc.edu.cn](mailto:{zkcys001,kevinh}@mail.ustc.edu.cn),  
[zhaodeli@gmail.com](mailto:zhaodeli@gmail.com), [jordan@cs.berkeley.edu](mailto:jordan@cs.berkeley.edu), [zhazj@ustc.edu.cn](mailto:zhazj@ustc.edu.cn).

## Abstract

The rank of neural networks measures information flowing across layers. It is an instance of a key structural condition that applies across broad domains of machine learning. In particular, the assumption of low-rank feature representations leads to algorithmic developments in many architectures. For neural networks, however, the intrinsic mechanism that yields low-rank structures remains vague and unclear. To fill this gap, we perform a rigorous study on the behavior of network rank, focusing particularly on the notion of rank deficiency. We theoretically establish a universal monotonic decreasing property of network rank from the basic rules of differential and algebraic composition, and uncover rank deficiency of network blocks and deep function coupling. By virtue of our numerical tools, we provide the first empirical analysis of the per-layer behavior of network rank in practical settings, *i.e.*, ResNets, deep MLPs, and Transformers on ImageNet. These empirical results are in direct accord with our theory. Furthermore, we reveal a novel phenomenon of independence deficit caused by the rank deficiency of deep networks, where classification confidence of a given category can be linearly decided by the confidence of a handful of other categories. The theoretical results of this work, together with the empirical findings, may advance understanding of the inherent principles of deep neural networks.

## 1 Introduction

In mathematics, the rank of a smooth function measures the volume of independent information captured by the function [21]. Deep neural networks are highly smooth functions, thus the rank of a network has long been an essential concept in machine learning that underlies many tasks such as information compression [48, 56, 36, 54, 49], network pruning [32, 55, 5, 25, 9], data mining [6, 24, 10, 57, 18, 29], computer vision [59, 58, 31, 27, 29, 60], and natural language processing [8, 28, 7, 11]. Numerous methods are either designed to utilize the mathematical property of network ranks, or are derived from an assumption that low-rank structures are to be preferred.

Yet a rigorous investigation to the behavior of rank of general networks, combining both theoretical and empirical arguments, is still absent in current research, weakening our confidence in the being able to predict performance. To the best of our knowledge, there are only a few previous works discussing the rank behavior of specific network architectures, like attention blocks [14] and BatchNorms [12, 4] in pure MLP structures. The empirical validation of those methods are also limited to shallow networks, specific architectures, or merely the final layers of deep networks, leaving the global behavior of general deep neural networks mysterious due to prohibitive space-time complexity for measuring them. Rigorous work on network rank that combines both strong theoretical and empirical evidence would have significant implications.

In this paper, we make several contributions towards this challenging goal. We find that the two essential ingredients of deep learning, the chain rules of differential operators and matrix multiplications, are enough to establish a universal principle—that network rank decreases monotonically with the depth of networks. Two factors further enhance the speed of decreasing: a) the explicit rank deficiency of many frequently used network modules, and b) an intrinsic potential of spectrum centralization enforced by the nature of coupling of massive composite functions. To empirically validate our theory, we design numerical tools to efficiently and economically examine the rank behavior of deep neural networks. This is a non-trivial task, as rank is very sensitive to noise and perturbation, and computing ranks of large networks is computationally prohibitive in time and space. Finally, we uncover an interesting phenomenon of independence deficit in multi-class classification networks. We find that many classes do not have their own unique representations in the classification network, and some highly irrelevant classes can decide the outputs of others. This independence deficit can significantly deteriorate the performance of networks in generalized data domains where each class demands a unique representation. In conclusion, the results of this work, together with the numerical tools we invent, may advance understanding of intrinsic properties of deep neural networks, and provide foundations for a broad study of low-dimensional structures in machine learning.

## 2 Preliminaries

**Settings** We consider the general deep neural network with  $L$  layers. It is a smooth vector-valued function  $\mathbf{F} : \mathbb{R}^n \rightarrow \mathbb{R}^d$ , where  $\mathbb{R}^n$  and  $\mathbb{R}^d$  are the ambient space of inputs and outputs, respectively. Deep neural networks are coupling of multiple layers, thus we write  $\mathbf{F}$  as:

$$\mathbf{F} = \mathbf{f}^L \circ \mathbf{f}^{L-1} \circ \cdots \circ \mathbf{f}^1. \quad (1)$$

For simplicity, we further write the  $k$ -th sub-network<sup>1</sup> of  $\mathbf{F}$  as

$$\mathbf{F}_k = \mathbf{f}^k \circ \cdots \circ \mathbf{f}^1, \quad (2)$$

and we use  $\mathcal{F}_k = \mathbf{F}_k(\mathcal{X})$  to denote the feature space of the  $k$ -th sub-network on the data domain  $\mathcal{X}$ . We are more interested in the behavior of network rank in the feature spaces rather than scalar outputs (which trivially have rank 1). Thus for classification or regression networks that output a scalar value, we will consider  $\mathbf{F} = \mathbf{F}_L$  as the transformation from the input space to the final feature space instead. Thus, we always have  $n \gg 1$  and  $d \gg 1$ . For example, for ResNet-50 [19] architecture on ImageNet, we only consider the network slice from the inputs to the last feature layer of 2,048 units.

**Rank of Function** The rank of a function  $\mathbf{f} = (\mathbf{f}_1, \dots, \mathbf{f}_d)^T : \mathbb{R}^n \rightarrow \mathbb{R}^d$  refers to the rank of its Jacobian matrix  $J_{\mathbf{f}}$  over its input domain  $\mathcal{X}$ , which is defined as

$$\text{Rank}(\mathbf{f}) = \text{Rank}(J_{\mathbf{f}}) = \text{Rank}((\partial \mathbf{f}_i(\mathbf{x}) / \partial x_j)_{n \times d}). \quad (3)$$

The rank of a function represents the volume of information captured by it in the output [21]. That is why it is so important to investigate the behavior of neural networks and many practical applications. Theoretically, by the rank theorem and Sard's theorem of manifolds [21], we can know that rank of the function equals the intrinsic dimension of its output feature space, as captured by the following lemma.<sup>2</sup>

**Lemma 1.** Suppose that  $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^d$  is smooth almost everywhere. Let  $\text{Rank}(\mathbf{f}) = r$ . If data domain  $\mathcal{X}$  is a manifold embedded in  $\mathbb{R}^n$  and  $\phi : \mathcal{U} \rightarrow \mathcal{O}$  is a smooth bijective parameterization from an open subset  $\mathcal{U} \subset \mathbb{R}^s$  to an open subset  $\mathcal{O} \subset \mathcal{X}$ , then we have  $\dim(\mathbf{f}(\mathcal{X})) = \text{Rank}(J_{\mathbf{f} \circ \phi}) \leq r$ . Thus, the rank of function  $\mathbf{f}$  gives an upper bound for the intrinsic dimension  $\dim(\mathbf{f}(\mathcal{X}))$  of the output space.

It is worth mentioning that the intrinsic dimension  $\dim(\mathbf{f}(\mathcal{X}))$  of the feature space is usually hard to measure, so the rank of the network gives an operational estimate of it.

---

<sup>1</sup>In this paper, sub-network means network slice from the input to some intermediate feature layer; layer network means an independent component of the network, without skip connections from the outside to it, like bottleneck layer of ResNet-50.

<sup>2</sup>Due to space limitation, all the related proofs are attached in the Appendix.

### 3 Numerical Tools

Validating the rank behavior of deep neural networks is a challenging task because it involves operations of high complexity on large-scale non-sparse matrices, which is infeasible both in time and space. Computing the full Jacobian representation of sub-networks of ResNet-50, for example, consumes over 150G GPU memory and several days at a single input point. In accuracy, this is even more challenging as rank is very sensitive to small perturbations. The numerical accuracy of float32,  $1.19e - 7$  [40], cannot be trivially neglected in computing matrix ranks. Thus, in this section we establish some numerical tools for validating our subsequent arguments, and provide rigorous theoretical support for them.

#### 3.1 Numerical Rank: Stable Alternative to Rank

The rank of large matrices is known to be unstable: it varies significantly under even small noise perturbations [42]. Matrices perturbed by even small Gaussian noises are almost surely of full rank, regardless of the true rank of the original matrix. Thus in practice we have to use an alternative: we count the number of singular values larger than some given threshold  $\epsilon$  as the numerical rank of the matrix. Let  $\mathbf{W} \in \mathbb{R}^{n \times d}$  be a given matrix. Its numerical rank with tolerance  $\epsilon$  is

$$\text{Rank}_\epsilon(\mathbf{W}) = \#\{i \in \mathbb{N}_+ : i \leq \min\{n, d\}, \sigma_i \geq \epsilon \|\mathbf{W}\|_2\}, \quad (4)$$

where  $\|\mathbf{W}\|_2$  is the  $\ell_2$  norm (spectral norm) of matrix  $\mathbf{W}$ ,  $\sigma_i, i = 1, \dots, \min\{n, d\}$  are its singular values, and  $\#$  is the counting measurement for finite sets. We can prove that the numerical rank is stable under small perturbations. Based on Weyl inequalities [50], we have the following theorem.

**Theorem 1.** *For any given matrix  $\mathbf{W}$ , almost every tolerance  $\epsilon > 0$ , and any perturbation matrix  $\mathbf{D}$ , there exists a positive constant  $\delta_{\max}(\epsilon)$  such that  $\forall \delta \in [0, \delta_{\max}(\epsilon)]$ ,  $\text{Rank}_\epsilon(\mathbf{W} + \delta \mathbf{D}) = \text{Rank}_\epsilon(\mathbf{W})$ . If  $\mathbf{W}$  is a low-rank matrix without random perturbations, then there is a  $\epsilon_{\max}$  such that for any  $\epsilon < \epsilon_{\max}$ ,  $\text{Rank}_\epsilon(\mathbf{W} + \delta \mathbf{D}) = \text{Rank}_\epsilon(\mathbf{W}) = \text{Rank}(\mathbf{W})$  for all  $\delta \in [0, \delta_{\max}(\epsilon)]$ .*

This property of the numerical rank metric makes it a suitable tool for investigating the rank behavior of neural networks. Possible small noises can be filtered out in Jacobian matrices of networks by using numerical rank. It is worth mentioning that random matrices no longer have full rank almost surely under the numerical rank. Instead their rank distribution can be inferred from the well-known Marcenko–Pastur distribution [34] of random matrices. So under numerical rank, low-rank matrices will be commonly seen. In this paper, we always use the numerical rank when measuring ranks.

#### 3.2 Partial Rank of the Jacobian: Estimating Lower Bound of Lost Rank in Deep Networks

To enable the validation of trend of the network ranks, we propose to compute only the rank of sub-matrices of the Jacobian as an alternative. Those sub-matrices are also the Jacobian matrices with respect to a fixed small patch of inputs. Rigorously, given a function  $\mathbf{f}$  and its Jacobian  $\mathbf{J}_f$ , we denote partial rank of the Jacobian as the rank of a sub-matrix of the Jacobian that consists of the  $j_1$ -th,  $j_2$ -th, ...,  $j_K$ -th column of the original Jacobian

$$\text{PartialRank}(\mathbf{J}_f) = \text{Rank}(\text{Sub}(\mathbf{J}_f, j_1, \dots, j_K)) = \text{Rank}((\partial f_i / \partial x_{j_k})_{d \times K}), \quad (5)$$

where  $1 \leq j_1 < \dots < j_K \leq n$ . We can efficiently compute sub-matrix of the Jacobian by zero padding to small patches of input images. For any data point  $\mathbf{x} \in \mathbb{R}^n$ , let  $\text{Sub}(\mathbf{x}, j_1, \dots, j_K) = (\mathbf{x}_{j_1}, \dots, \mathbf{x}_{j_K})^T \in \mathbb{R}^K$ , and  $\psi$  pad  $\text{Sub}(\mathbf{x}, j_1, \dots, j_K)$  to the spatial size of  $\mathbf{x}$  with zeros:  $\psi(\text{Sub}(\mathbf{x}, j_1, \dots, j_K)) = (0, \dots, 0, \mathbf{x}_{j_1}, 0, \dots, \mathbf{x}_{j_K}, 0, \dots, 0)^T \in \mathbb{R}^n$  with  $\psi(\text{Sub}(\mathbf{x}, j_1, \dots, j_K))_{j_k} = \mathbf{x}_{j_k}, k = 1, \dots, K$ . We then have  $\mathbf{J}_{f \circ \psi} = \text{Sub}(\mathbf{J}_f, j_1, \dots, j_K)$ . As  $K$  can be very small compared with  $n$ , computing  $\mathbf{J}_{f \circ \psi}$  can be very cheap in time and space. The partial rank of Jacobian matrices of the network layers measures information captured among the spatial footprint  $j_1, \dots, j_K$  of the original input. They inherit the order relation of the rank of full Jacobian matrices. Thus we can validate the rank diminishing of network Jacobian matrices through the partial rank.

**Lemma 2.** *For differentiable  $f_1, f_2$ ,  $|\text{Rank}(f_1) - \text{Rank}(f_2 \circ f_1)| \geq |\text{Rank}(\text{Sub}(f_1, j_1, \dots, j_K)) - \text{Rank}(\text{Sub}(f_2 \circ f_1, j_1, \dots, j_K))|, \forall 1 \leq K \leq n, 1 \leq j_1, \dots, j_K \leq n$ . Thus variance of partial ranks of adjacent sub-networks gives a lower bound on the variance of their ranks.*

### 3.3 Classification Dimension: Estimating Final Feature Dimension

Measuring the intrinsic dimension of feature manifolds is known to be intractable. So we turn to an approximation procedure. For most classification networks, a linear regression over the final feature manifold decides the final network prediction and accuracy. So we can estimate the intrinsic dimension as the minimum number of principal components in the final feature space to preserve a high classification accuracy. Given network slice  $\mathbf{F} : \mathbb{R}^n \rightarrow \mathbb{R}^d$  from input  $\mathcal{X} \subset \mathbb{R}^n$  to final feature space  $\mathbf{F}(\mathcal{X}) \subset \mathbb{R}^d$ , we independently sample  $N$  points from random variable  $\mathbf{F}(\mathbf{x})$ ,  $\mathbf{x} \sim \mathbb{P}_{\mathcal{X}}$ , where  $\mathbb{P}_{\mathcal{X}}$  is the distribution of validation set of data  $\mathcal{X}$ . We then compute the covariance matrix  $\Sigma$  of those  $N$  samples, and eigenvectors  $\mathbf{q}^1, \dots, \mathbf{q}^d$  of  $\Sigma$ , sorted by their eigenvalues  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d$ . Let  $\text{cls} : \mathbb{R}^d \rightarrow \mathbb{R}^c$  be the classification predictions based on the final feature representation  $\mathbf{F}(\mathbf{x})$ ,  $\text{Pro}_k$  be projection operator in Euclidean space that projects to the linear subspace spanned by top- $K$  eigenvectors  $\mathbf{q}^1, \dots, \mathbf{q}^k$ ,  $k \leq d$ ,  $\mathbb{P}_{\mathbf{x}, \mathbf{y}}$  be the joint distribution of sample  $\mathbf{x}$  and its label  $\mathbf{y}$ , and  $\mathbf{1}_{\text{cond}}$  the indicator for condition cond. The classification dimension is then defined as

$$\text{ClsDim}(\mathbf{F}(\mathcal{X})) = \min_k \{k : \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \mathbb{P}_{\mathcal{X}, \mathcal{Y}}} [\mathbf{1}_{\text{cls}(\text{Pro}_k(\mathbf{F}(\mathbf{x}))) == \mathbf{y}}] \geq 1 - \epsilon\}, \quad (6)$$

which is the minimum dimensionality needed to reconstruct the classification accuracy of the whole model.

## 4 Principle of Rank Diminishing

We turn to the *principle of rank diminishing*. We first give a universal justification with minimum limitation on the network, so that we can safely apply this principle to many practical scenarios.

The principle of rank diminishing describes the behavior of general neural networks with almost everywhere smooth components, which exhibits the monotonic decreasing of network ranks and intrinsic dimensionality of feature manifolds as follows.

**Theorem 2** (Principle of Rank Diminishing). *Suppose that each layer  $f_i$ ,  $i = 1, \dots, L$  of network  $\mathbf{F}$  is almost everywhere smooth and data domain  $\mathcal{X}$  is a manifold, then both the rank of sub-networks and intrinsic dimension of feature manifolds decrease monotonically by depth:*

$$\text{Rank}(f_1) \geq \text{Rank}(f_2 \circ f_1) \geq \dots \geq \text{Rank}(f_{L-1} \circ \dots \circ f_1) \geq \text{Rank}(\mathbf{F}_L), \quad (7)$$

$$\dim(\mathcal{X}) \geq \dim(\mathcal{F}_1) \geq \dim(\mathcal{F}_2) \geq \dots \geq \dim(\mathcal{F}_L). \quad (8)$$

**Short Argument that the Principle Should Hold Universally.** Theorem 2 is ultra intrinsic for deep neural networks. It comes directly from the chain rules of differential operators and basic rules of matrix multiplications. The basic rule of matrix multiplication tells that, for any two matrices  $\mathbf{A}$  and  $\mathbf{B}$ , we have  $\text{Rank}(\mathbf{AB}) \leq \min\{\text{Rank}(\mathbf{A}), \text{Rank}(\mathbf{B})\}$  [22]. Taking this into the chain rule of differential of  $\mathbf{J}_F = \mathbf{J}_{f^L} \mathbf{J}_{f^{L-1}} \dots \mathbf{J}_{f^1}$ , we then have  $\text{Rank}(\mathbf{J}_{F_k}) = \text{Rank}(\mathbf{J}_{f^k \circ F_{k-1}}) = \text{Rank}(\mathbf{J}_{f^k} \mathbf{J}_{F_{k-1}}) \leq \text{Rank}(\mathbf{J}_{F_{k-1}})$ ,  $k = 2, \dots, L$ , which is Eq. (7). Applying Lemma 1 to Eq. (7) then yields Eq. (8).

**Chance of Equal Sign Holding is Small.** A hypothetical but not practical concern would be that, is it possible that most of the equal signs of Eqs. (7) and (8) hold, so that the rank of network remains no significant dropping throughout the network? This concern can be mitigated by empirical and theoretical arguments. In what follows we will find that, 1) in practice, the rank of sub-networks decreases significantly after applying subsequent layers as shown in Fig. 1, and 2) in theory, there are two strong impetuses in deep neural networks to enforce the strict decreasing of ranks which we will discuss in Secs. 4.1 and 4.2.

### 4.1 Structural Impetus of Strict Decreasing

Numerous explicit structures of the network layers can lead to a strict decrease in network ranks. Specifically, the following theorem gives a condition for the strictly greater signs to hold in the principle of rank diminishing.

**Theorem 3.** <sup>3</sup> Roughly speaking, if almost everywhere on the input feature manifold, there is a direction such that moving along this direction keeps the output invariant, then the intrinsic dimension

---

<sup>3</sup>The rigorous version is given in the Appendix.

Arch.	Network	Activ.	#Param.	Main Block	#Layer	Top-1 Acc.
ResNets	ResNet-18 [19]	ReLU [37]	11.7M	Bottleneck	11	69.8%
	ResNet-50 [19]	ReLU [37]	25.6M	Bottleneck	19	76.1%
MLP-like	GluMixer-24 [43]	SiLU [20]	25.0M	Mixer-Block	24	78.1%
	ResMLP-S24 [46]	GELU [20]	30.0M	Mixer-Block	24	79.4%
Transformer	ViT-T [15]	GELU [20]	5.7M	ViT-Block	13	75.5%
	Swin-T [33]	GELU [20]	29.0M	Swin-Block	18	81.3%

Table 1: Information of networks used in empirical validations. All pretrained on ImageNet.

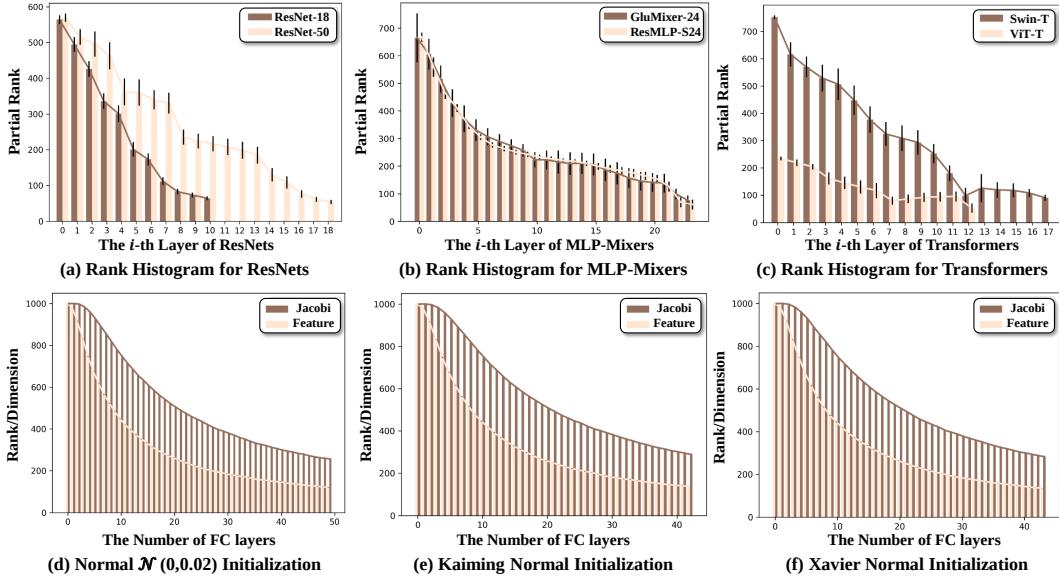


Figure 1: Partial rank of Jacobian matrices of CNN, MLP, and Transformer architecture networks for different layers on ImageNet (top row); rank of Jacobian matrices and feature dimensions of linear MLP network following conditions of Theorem 5 (bottom rule). All the models show a similar trend of exponential decreasing of ranks as predicted by Theorems 4 and 5.

of the output feature manifold will be strictly lower than that of the input. The maximum number of independent such directions gives a lower bound on the number of lost intrinsic dimensions.

By this theorem, one can immediately find that most frequently used layer designs have high risk in inducing the strict decreasing of network ranks. Normalization layers like LayerNorm [2], InstanceNorm [47], and BatchNorm [26] may lose dimensions modestly, as the output feature remains invariant along the normalized direction at each point. Linear layers like convolutions, linear transformations (e.g. dense layers), and attentions, can lose rank considerably according to the rank of their weight matrices. They constitute the explicit structural impetus to decrease network ranks and intrinsic dimensions of feature manifolds.

## 4.2 Implicit Impetus of Strict Decreasing

Apart from the structural impetus we propose in Sec. 4.1, there is a more intrinsic strength to pull down network ranks, which we call the implicit impetus. Deep neural networks repeatedly apply layer networks from a fixed function pool (ReLU, MLP, CNN, attention, ResNet block, etc.) to the input data and intermediate features to get outputs. Such paradigm accords with the cocycle dynamic systems studied by Lyapunov *et al.* [44, 52], where the Furstenberg–Kesten theorem [16] and multiplicative ergodic theorem [41] prove that logarithms of singular values divided by evolution time of such chaos system converge to stable constants when time goes to infinity. While products of long chains of matrices are the simplest form of cocycle dynamic systems [30], we can get an intrinsic impetus of rank collapse tendency of Jacobian matrices independent of network architectures.

**Theorem 4** (Spectrum Centralization of Function Coupling). *Let the network be  $F = f^L \circ \dots \circ f^1$ , and all the ambient dimensions of feature manifolds be the same as the ambient dimension of inputs,*

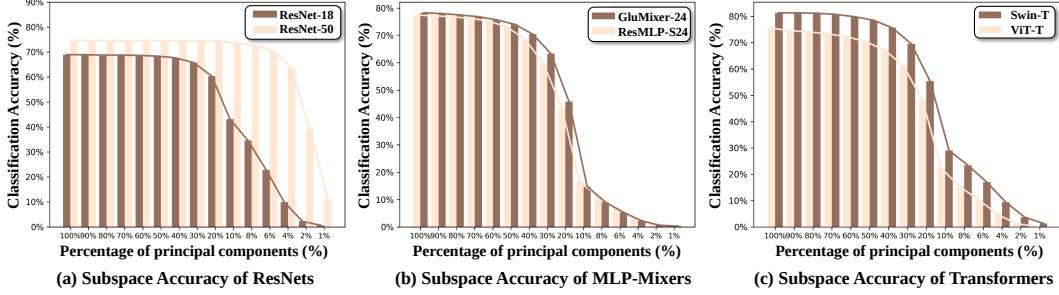


Figure 2: Classification Accuracy (top-1) of using subspaces spanned by top- $k\%$  eigenvectors (principal components) of the final feature manifolds. For all networks a small percentage (see Tab. 2) of eigenvectors are enough to reproduce the classification accuracy of the whole network, indicating a low intrinsic dimension of final feature manifolds. **Note that the  $x$ -axes are non-linear.**

Networks	ResNet-18	ResNet-50	GluMixer-24	ResMLP-S24	Swin-T	ViT-T
ClsDim	149	131	199	196	344	109
Ambient Dim.	512	2048	384	384	768	192

Table 2: Classification dimensions (with respect to 95% classification performance of the ambient feature space  $\mathbb{R}^d$ ) and ambient dimensions of the final feature manifolds of different networks. All networks have low intrinsic dimensions for final features.

i.e.,  $f^k : \mathbb{R}^n \rightarrow \mathbb{R}^n, k = 1, \dots, L$ . Suppose the Jacobian matrix of each layer  $f_i$  independently follows some distribution  $\mu$ , and  $\mathbb{E}_\mu[\max\{\log \|J_{f^k}^{\pm 1}\|_2, 0\}] < \infty$ . Let  $\sigma_k$  denote the  $k$ -th largest singular value of  $J_F$ . Then there is an integer  $r < n$  and positive constants  $\mu_r, \dots, \mu_n$  that only depend on  $\mu$  such that we have for  $\mu$ -almost everywhere,

$$\frac{\sigma_k}{\|J_F\|_2} \sim \exp(-L\mu_k) \rightarrow 0, k = r, \dots, n, \text{ as } L \rightarrow \infty, \quad (9)$$

meaning that for any tolerance  $\epsilon > 0$ ,  $\text{Rank}_\epsilon(F)$  drops below  $r + 1$  with an exponential speed as  $L \rightarrow \infty$ .

If further assuming that the elements of Jacobian matrices follow Gaussian distributions, we can prove  $r = 1$  and give a more accurate estimation of constants  $\mu_r, \dots, \mu_n$ . As a consequence, we can find that rank of networks collapses to 1 almost surely, which is formalized in the following theorem.

**Theorem 5.** Let the network be  $F = f^L \circ \dots \circ f^1$ , and all the ambient dimensions of feature manifolds be the same as the ambient dimension of inputs, i.e.,  $f^k : \mathbb{R}^n \rightarrow \mathbb{R}^n, k = 1, \dots, L$ . Suppose that  $J_{f^k}$  independently follows the standard Gaussian distribution. Let  $\sigma_k$  denote the  $k$ -th largest singular value of  $J_F$ . Then almost surely

$$\lim_{L \rightarrow \infty} \left( \frac{\sigma_k}{\|J_F\|_2} \right)^{\frac{1}{L}} = \exp \frac{1}{2} \left( \psi \left( \frac{n-k+1}{2} \right) - \psi \left( \frac{n}{2} \right) \right) < 1, k = 2, \dots, n, \quad (10)$$

where  $\psi = \Gamma/\Gamma'$  and  $\Gamma$  is the Gamma function. That means for a large  $L$  and any tolerance  $\epsilon$ ,  $\text{Rank}_\epsilon(F)$  drops to 1 exponentially with speed  $nC^L$ , where  $C < 1$  is a positive constant that only depends on  $n$ .

**Connection with Gradient Explosion** Bengio *et. al.* [3, 39] discuss the gradient explosion issue of deep neural networks, where the largest singular value of the Jacobian matrix tends to infinity when the layer gets deeper. This problem could be viewed as a special case of Theorem 5 that investigates the behavior of all singular values of deep neural networks. The behavior of network ranks in fact manipulates the well-known gradient explosion issue. Rigorously, we have the following conclusion.

**Corollary 1.** Under the condition of Theorem 5, then almost surely gradient explosion happens at an exponential speed, i.e.,  $\log \|J_F\|_2 = \log \sigma_1 \sim \frac{L}{2}(\log 2 + \psi(n/2)) \rightarrow \infty$  when  $L$  is large.

### 4.3 Validation

**Setup** In this section, we validate our theory in three types of architectures of benchmark deep neural networks, CNNs, MLPs, and Transformers, in the ImageNet [13] data domain. Information of

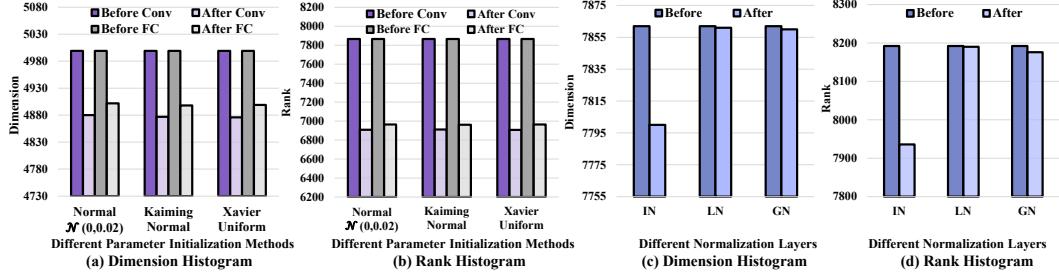


Figure 3: PCA dimension of feature spaces and rank of Jacobian matrix for commonly seen network components under standard Gaussian inputs and randomized weights. Convolution and FC layers tend to lose rank considerably; normalization layers, like InstanceNorm (IN) [47], LayerNorm (LN) [2], and GroupNorm (GN) [53], lose rank modestly. But none can preserve rank.

those networks is listed in Tab. 1. For validating the tendency of network rank of Jacobian matrices, we use the numerical rank of sub-matrices of Jacobian on the central  $16 \times 16 \times 3$  image patch of input images. We report the results of other choices of patches in the Appendix. When measuring rank, we set  $\epsilon = \text{eps} \times N$ , where  $\text{eps}$  is the digital accuracy of float32 (*i.e.*,  $1.19e - 7$ ) and  $N$  is the number of singular values of the matrix to measure. This threshold represents the minimum digital accuracy of numerical rank we can capture in data stored as float32. All the experiments are conducted on the validation set of ImageNet and NVIDIA A100-SXM-80G GPUs.

**Rank Diminishing of Jacobians** As is discussed in Sec. 3.2 and Lemma 2, the partial rank of the Jacobian is a powerful weapon for us to detect the behavior of huge Jacobian matrices, which are infeasible to compute in practice. The decent value of partial ranks of adjacent sub-networks provides a lower bound to that of full ranks of them. Fig. 1 (a,b,c) report the partial rank of Jacobian matrices of three types of architectures, where we can find consistent diminishing of partial ranks in each layer, indicating a larger rank losing for the full rank of Jacobian matrices.

**Intrinsic Dimension of the Final Feature Manifold** To get a further estimation of how many dimensions remain in the final feature representation, we measure the classification dimension in Fig. 2 and Tab. 2. We report the classification accuracy produced by projecting final feature representations to its top  $k\%$  eigenvectors in Fig. 2. We choose a threshold of  $\epsilon$  such that this procedure can reproduce 95% of the original accuracy of the network. The corresponding ClsDim is reported in Tab. 2. As discussed in Sec. 3.3, this gives an estimation of the intrinsic dimension of the final feature manifold. We can find a universal low-rank structure for all types of networks.

**Implicit Impetus** Theorem 5 gives an exponential speed of rank decent by layers. We find that it corresponds well with practice. We investigate this exponential law in a toy network of MLP-50, which is composed of 50 dense layers, each with 1,000 hidden units. The MLP-50 network takes Gaussian noise vectors of  $\mathbb{R}^{1000}$  as inputs, and returns a prediction of 1,000 categories. As all the feature manifolds are linear subspaces in this case, their intrinsic dimensions can be directly measured by the numerical rank of their covariance matrices. We report the full rank of Jacobian matrices and intrinsic dimensions of feature manifolds under three different randomly chosen weights in Fig. 1 (d,e,f). Due to the digital accuracy of float32, we stop calculation in each setting when the absolute values of elements of the matrices are lower than  $1.19e - 7$ . We can find standard curves of exponential laws in all cases for both ranks of Jacobian and intrinsic dimensions of features. By comparison, we can further find that the ranks of benchmark deep neural networks on ImageNet bear a striking resemblance to the exponential laws of our toy setting, which confirms the proposed implicit impetus in those models.

**Structural Impetus** We validate the structural impetus in Fig. 3. To give an estimation for general cases, here we use Gaussian noises with the size of  $128 \times 8 \times 8$  as inputs, and randomize weights of the network components to be validated. We plug those components into a simple fully-connected (FC) layer of 8,192 hidden units. As the structure is simple, we directly measure the intrinsic dimension of feature spaces and the full rank of Jacobian matrices before and after the features pass the network components to be measured. The dimension is determined by the number of PCA eigenvalues [23, 51] larger than  $1.19e - 7 \times N \times \sigma_{\max}$ , where  $N$  is the number of PCA eigenvalues, and  $\sigma_{\max}$  is the

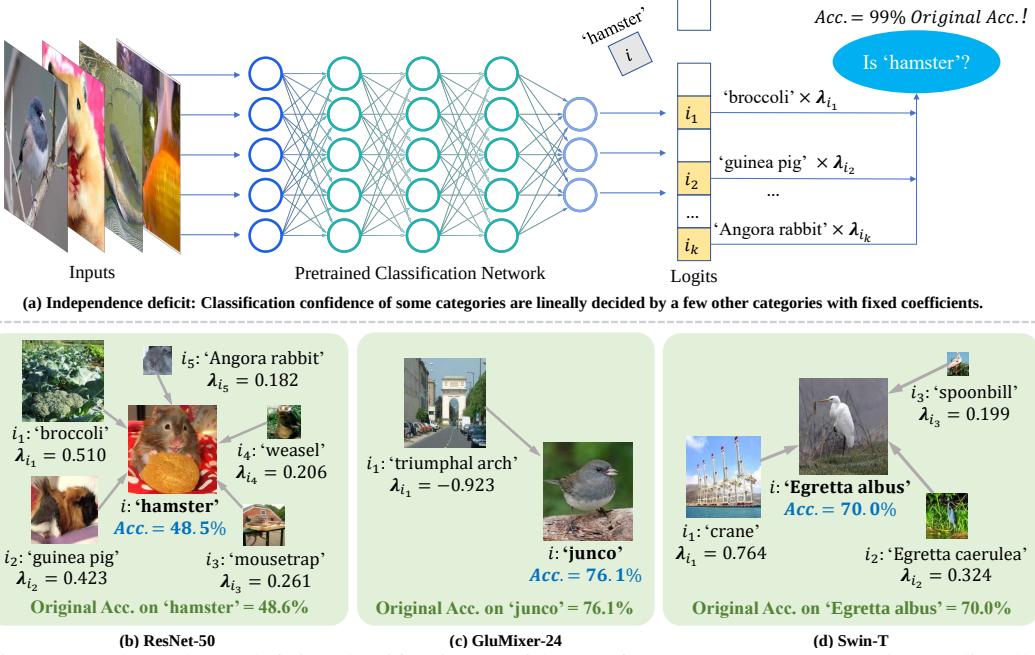


Figure 4: Independence deficit. Classification confidence of some ImageNet categories are linearly decided by a few other categories with fixed coefficients in the whole data domain. We illustrate this phenomenon in (a). Here we present some results from ResNet-50, GluMixer-24, and Swin-T. In the (b,c,d) we illustrate the categories of  $i_1, \dots, i_k$  (in the surrounding) to linearly decide category  $i$  (in the center) and their corresponding weights  $\lambda_{i_1}, \dots, \lambda_{i_k}$ . The classification accuracy on the validation set of using Eq. (12), instead of the true logits, to predict the label is reported in blue (if tested on positive samples only, the accuracy rates are 98%, 90%, 82% for cases in (b,c,d) correspondingly). For comparison, the original accuracy for the corresponding categories are reported in green. We can find that 1) a few other categories can decide the confidence of the target category  $i$ ; 2) some very irrelevant categories contribute the largest weights. For example in (c), the logits of class ‘junco’ is the negative of ‘triumphal arch’. Both of them indicate a rather drastic competition of different categories for independent representations in final features due to the tight rank budgets.

largest PCA eigenvalue. The batch size is set to 5,000. We find that the convolution (the kernel size is  $3 \times 3$ ) and FC layers (the weight size is 8,192) tend to lose rank considerably, while different normalization layers also lose rank modestly. But none of them can preserve rank invariant.

**Possible Remission Approaches to Rank Diminishing** There are quite some techniques, at least in theory, can remiss the network rank diminishing. Typical examples are skip connection [14] and BatchNorm [12], which we will discuss in the Appendix due to page limitation.

## 5 Independence Deficit of Final Feature Manifolds

In this section, we provide a further perspective to study the low-rank structure of the final feature manifold, which induces an interesting finding of independence deficit in deep neural networks. We have already known that the final feature representations of deep neural networks admit a very low intrinsic dimension. Thus there are only a few independent representations to decide the classification scores for all the 1,000 categories of ImageNet. It is then curious whether we can predict the outputs of the network for some categories based on the outputs for a few other categories, as illustrated in Fig. 4 (a). And if we can, will those categories be strongly connected to each other? A surprising fact is that, we can find many counter examples of irrelevant categories dominating the network outputs for given categories regarding various network architectures. This interesting phenomenon indicates a rather drastic competing in the final feature layer for the tight rank budgets of all categories, which yields non-realistic dependencies of different categories.

To find the dependencies of categories in final features, we can solve the following Lasso problem [45],

$$\boldsymbol{\lambda}^* = \arg \min_{\boldsymbol{\lambda}_{i=-1}} \mathbb{E}_{\mathbf{x}} [\|\boldsymbol{\lambda}^T \mathbf{W} \mathbf{F}(\mathbf{x})\|_2^2] + \eta \|\boldsymbol{\lambda}\|_1, \quad (11)$$

where  $\mathbf{F}(\mathbf{x}) \in \mathbb{R}^{1000}$  is the slice of network from inputs to the final feature representation,  $\mathbf{x}$  is the sample from ImageNet  $\mathcal{X}$ , and  $\mathbf{W}$  is the final dense layer. The solution  $\boldsymbol{\lambda}^*$  will be a sparse vector, with  $k$  non-zero elements  $\lambda_{i_1} \geq \lambda_{i_2} \geq \dots \geq \lambda_{i_k}, k \ll 1000$ . We can then get

$$\text{logits}(\mathbf{x}, i) \approx \lambda_{i_1} \text{logits}(\mathbf{x}, i_1) + \dots + \lambda_{i_k} \text{logits}(\mathbf{x}, i_k), i \notin \{i_1, \dots, i_k\}, k \ll 1000, \forall \mathbf{x} \in \mathcal{X}, \quad (12)$$

where  $\text{logits}(\mathbf{x}, i_j), j = 1, \dots, k$  is the logits of network for category  $i_j$ , i.e.,  $\text{logits}(\mathbf{x}, i_j) = \mathbf{W}_{i_j} \mathbf{F}(\mathbf{x})$ . It is easy to see that outputs for category  $i$  are linearly decided by outputs for  $i_1, \dots, i_k$  and are dominated by outputs for  $i_1$ .

In Fig. 4 we demonstrate the solutions of Eq. (12) for three different categories in ImageNet with  $\eta = 20$ , and network architectures ResNet-50, GluMixer-24, and Swin-T. The results are surprising. It shows that many categories of the network predictions are in fact ‘redundant’, as they are purely decided by the predictions of the other categories with simple linear coefficients. In this case, the entanglement of different categories cannot be avoided, thus the network may perform poorly under domain shift. An even more surprising finding is that, some very irrelevant categories hold the largest weights when deciding the predictions of the redundant categories, which means that the networks just neglect the unique representations of those categories in training and yield over-fitting when predicting them.

## 6 Related Work

Previous studies of rank deficiency in deep neural networks follow two parallel clues. One is the study of rank behavior in specific neural network architectures. [14] studies deep networks consisting of pure self-attention networks, and proves that they converge exponentially to a rank-1 matrix under the assumption of globally bounded weight matrices. [12] studies the effect of BatchNorm on MLPs and shows that BatchNorm can prevent drastic diminishing of network ranks in some small networks and datasets. Both of those works avoid directly validating the behavior of network ranks in intermediate layers due to the lacking of efficient numerical tools. An independent clue is the study of implicit self-regularization, which finds that weight matrices tend to lose ranks after training. [35] studies this phenomenon in infinitely-wide, over-parametric neural networks with tools from random matrix theory. [1] studies this phenomenon in deep matrix decomposition. Those works focus on the theoretical behavior of rank of weight matrices induced by the training instead of network ranks.

## 7 Conclusion

This paper studies the rank behavior of deep neural networks. In contrast to previous work, we focus on directly validating rank behavior with deep neural networks of diverse benchmarks and various settings for real scenarios. We first formalize the analysis and measurement of network ranks. Then under the proposed numerical tools and theoretical analysis, we demonstrate the universal rank diminishing of deep neural networks from both empirical and theoretical perspectives. We further support the rank-deficient structure of networks by revealing the independence deficit phenomenon, where network predictions for a category can be linearly decided by a few other, even irrelevant categories. The results of this work may advance understanding of the behavior of fundamental network architectures and provide intuition for a wide range of work pertaining to network ranks.

## References

- [1] S. Arora, N. Cohen, W. Hu, and Y. Luo. Implicit regularization in deep matrix factorization. *Adv. Neural Inform. Process. Syst.*, 32, 2019.
- [2] J. L. Ba, J. R. Kiros, and G. E. Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016.
- [3] Y. Bengio, P. Simard, and P. Frasconi. Learning long-term dependencies with gradient descent is difficult. *IEEE Trans. Neural Netw.*, 5(2):157–166, 1994.
- [4] N. Bjorck, C. P. Gomes, B. Selman, and K. Q. Weinberger. Understanding batch normalization. *Adv. Neural Inform. Process. Syst.*, 31, 2018.

- [5] C. Blakeney, Y. Yan, and Z. Zong. Is pruning compression?: Investigating pruning via network layer similarity. In *IEEE Winter Conf. Appl. Comput. Vis.*, pages 914–922, 2020.
- [6] X. Cai, C. Ding, F. Nie, and H. Huang. On the equivalent of low-rank linear regressions and linear discriminant analysis based regressions. In *Proc. ACM SIGKDD Int. Conf. Knowl. Discov. Data Min.*, pages 1124–1132, 2013.
- [7] B. Chen, T. Dao, E. Winsor, Z. Song, A. Rudra, and C. Ré. Scatterbrain: Unifying sparse and low-rank attention. *Adv. Neural Inform. Process. Syst.*, 34, 2021.
- [8] P. Chen, S. Si, Y. Li, C. Chelba, and C.-J. Hsieh. Groupreduce: Block-wise low-rank approximation for neural language model shrinking. *Adv. Neural Inform. Process. Syst.*, 31, 2018.
- [9] P. Chen, H.-F. Yu, I. Dhillon, and C.-J. Hsieh. Drone: Data-aware low-rank compression for large NLP models. *Adv. Neural Inform. Process. Syst.*, 34:29321–29334, 2021.
- [10] Y. Cheng, L. Yin, and Y. Yu. Lorslim: Low rank sparse linear methods for top-n recommendations. In *IEEE Int. Conf. on Data Min.*, pages 90–99. IEEE, 2014.
- [11] J. Chiu, Y. Deng, and A. Rush. Low-rank constraints for fast inference in structured models. *Adv. Neural Inform. Process. Syst.*, 34, 2021.
- [12] H. Daneshmand, J. Kohler, F. Bach, T. Hofmann, and A. Lucchi. Batch normalization provably avoids ranks collapse for randomly initialised deep networks. *Adv. Neural Inform. Process. Syst.*, 33:18387–18398, 2020.
- [13] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. ImageNet: A large-scale hierarchical image database. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 248–255. Ieee, 2009.
- [14] Y. Dong, J.-B. Cordonnier, and A. Loukas. Attention is not all you need: Pure attention loses rank doubly exponentially with depth. In *Int. Conf. Mach. Learn.*, pages 2793–2803. PMLR, 2021.
- [15] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *Int. Conf. Learn. Represent.*, 2020.
- [16] H. Furstenberg and H. Kesten. Products of random matrices. *Ann. Math. Stat.*, 31(2):457–469, 1960.
- [17] H. Furstenberg and H. Kesten. Random matrix products and measures on projective spaces. *Israel J. Math.*, 46:12–32, 1983.
- [18] D. Goldfarb and Z. Qin. Robust low-rank tensor recovery: Models and algorithms. *SIAM J. Matrix Anal. Appl.*, 35(1):225–253, 2014.
- [19] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 770–778, 2016.
- [20] D. Hendrycks and K. Gimpel. Gaussian error linear units (GELUs). *arXiv preprint arXiv:1606.08415*, 2016.
- [21] M. W. Hirsch. *Differential topology*, volume 33. Springer Science & Business Media, 2012.
- [22] K. Hoffman. *Linear algebra*. Englewood Cliffs, NJ, Prentice-Hall, 1971.
- [23] H. Hotelling. Relations between two sets of variates. *Biometrika*, 28(3/4):321–377, 1936.
- [24] C.-J. Hsieh, K.-Y. Chiang, and I. S. Dhillon. Low rank modeling of signed networks. In *Proc. ACM SIGKDD Int. Conf. Knowl. Discov. Data Min.*, pages 507–515, 2012.
- [25] Y.-C. Hsu, T. Hua, S. Chang, Q. Lou, Y. Shen, and H. Jin. Language model compression with weighted low-rank factorization. In *Int. Conf. Learn. Represent.*, 2021.
- [26] S. Ioffe and C. Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *Int. Conf. Mach. Learn.*, pages 448–456. PMLR, 2015.
- [27] L. Jing, L. Yang, J. Yu, and M. K. Ng. Semi-supervised low-rank mapping learning for multi-label classification. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 1483–1491, 2015.
- [28] R. Karimi Mahabadi, J. Henderson, and S. Ruder. Compacter: Efficient low-rank hypercomplex adapter layers. *Adv. Neural Inform. Process. Syst.*, 34, 2021.

- [29] M. Kheirandishfard, F. Zohrizadeh, and F. Kamangar. Deep low-rank subspace clustering. In *IEEE Conf. Comput. Vis. Pattern Recog. Worksh.*, pages 864–865, 2020.
- [30] J. F. C. Kingman. Subadditive ergodic theory. *Ann. Probab.*, pages 883–899, 1973.
- [31] S. Kong and C. Fowlkes. Low-rank bilinear pooling for fine-grained classification. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 365–374, 2017.
- [32] M. Lin, R. Ji, Y. Wang, Y. Zhang, B. Zhang, Y. Tian, and L. Shao. Hrank: Filter pruning using high-rank feature map. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 1529–1538, 2020.
- [33] Z. Liu, Y. Lin, Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, and B. Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 10012–10022, 2021.
- [34] V. A. Marčenko and L. A. Pastur. Distribution of eigenvalues for some sets of random matrices. *Math. USSR Sb.*, 1(4):457, 1967.
- [35] C. H. Martin and M. W. Mahoney. Implicit self-regularization in deep neural networks: Evidence from random matrix theory and implications for learning. *J. Mach. Learn. Res.*, 22(165):1–73, 2021.
- [36] J. Mu, R. Xiong, X. Fan, D. Liu, F. Wu, and W. Gao. Graph-based non-convex low-rank regularization for image compression artifact reduction. *IEEE Trans. Image Process.*, 29:5374–5385, 2020.
- [37] V. Nair and G. E. Hinton. Rectified linear units improve restricted Boltzmann machines. In *Int. Conf. Mach. Learn.*, pages 807–814. PMLR, 2010.
- [38] C. M. Newman. The distribution of Lyapunov exponents: Exact results for random matrices. *Commun. Math. Phys.*, 103(1):121–126, 1986.
- [39] R. Pascanu, T. Mikolov, and Y. Bengio. On the difficulty of training recurrent neural networks. In *Int. Conf. Mach. Learn.*, pages 1310–1318. PMLR, 2013.
- [40] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Adv. Neural Inform. Process. Syst.*, 32, 2019.
- [41] Y. B. Pesin. Characteristic Lyapunov exponents and smooth ergodic theory. *Russ. Math. Surv.*, 32(4):55, 1977.
- [42] A. Quarteroni, R. Sacco, and F. Saleri. *Numerical mathematics*, volume 37. Springer Science & Business Media, 2010.
- [43] N. Shazeer. GLU variants improve transformer. *arXiv preprint arXiv:2002.05202*, 2020.
- [44] R. Temam. *Infinite-dimensional dynamical systems in mechanics and physics*, volume 68. Springer Science & Business Media, 2012.
- [45] R. Tibshirani. Regression shrinkage and selection via the lasso. *J. R. Stat. Soc. Series B Stat. Methodol.*, 58(1):267–288, 1996.
- [46] H. Touvron, P. Bojanowski, M. Caron, M. Cord, A. El-Nouby, E. Grave, G. Izacard, A. Joulin, G. Synnaeve, J. Verbeek, et al. ResMLP: Feedforward networks for image classification with data-efficient training. *arXiv preprint arXiv:2105.03404*, 2021.
- [47] D. Ulyanov, A. Vedaldi, and V. Lempitsky. Instance normalization: The missing ingredient for fast stylization. *arXiv preprint arXiv:1607.08022*, 2016.
- [48] T. Vogels, S. P. Karimireddy, and M. Jaggi. Practical low-rank communication compression in decentralized deep learning. *Adv. Neural Inform. Process. Syst.*, 33:14171–14181, 2020.
- [49] H. Wang and N. Ahuja. Rank-r approximation of tensors using image-as-matrix representation. In *IEEE Conf. Comput. Vis. Pattern Recog.*, volume 2, pages 346–353. IEEE, 2005.
- [50] H. Weyl. Das asymptotische verteilungsgesetz der eigenwerte linearer partieller differentialgleichungen. *Math. Ann.*, 71:441–479, 1912.
- [51] S. Wold, K. Esbensen, and P. Geladi. Principal component analysis. *Chemom. Intell. Lab. Syst.*, 2(1-3):37–52, 1987.

- [52] A. Wolf, J. B. Swift, H. L. Swinney, and J. A. Vastano. Determining Lyapunov exponents from a time series. *Physica D*, 16(3):285–317, 1985.
- [53] Y. Wu and K. He. Group normalization. In *Eur. Conf. Comput. Vis.*, pages 3–19, 2018.
- [54] J. Ye. Generalized low rank approximations of matrices. *Mach. Learn.*, 61(1):167–191, 2005.
- [55] X. Yu, T. Liu, X. Wang, and D. Tao. On compressing deep models by low rank and sparse decomposition. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 7370–7379, 2017.
- [56] Z. Zha, X. Yuan, B. Wen, J. Zhou, J. Zhang, and C. Zhu. From rank estimation to rank approximation: Rank residual constraint for image restoration. *IEEE Trans. Image Process.*, 29:3254–3269, 2019.
- [57] M. Zhan, S. Cao, B. Qian, S. Chang, and J. Wei. Low-rank sparse feature selection for patient similarity learning. In *IEEE Int. Conf. on Data Min.*, pages 1335–1340. IEEE, 2016.
- [58] T. Zhang, B. Ghanem, S. Liu, C. Xu, and N. Ahuja. Low-rank sparse coding for image classification. In *Int. Conf. Comput. Vis.*, pages 281–288, 2013.
- [59] Y. Zhang, Z. Jiang, and L. S. Davis. Learning structured low-rank representations for image classification. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 676–683, 2013.
- [60] H. Zou, T. Hastie, and R. Tibshirani. Sparse principal component analysis. *J. Comput. Graph. Stat.*, 15(2):265–286, 2006.

## Appendix

### Contents

<b>A Proofs</b>	<b>13</b>
A.1 Proof to Lemma 1	13
A.2 Proof to Theorem 1	13
A.3 Proof to Lemma 2	14
A.4 Proof to Theorem 2	14
A.5 Proof to Theorem 3	15
A.6 Proof to Theorem 4	15
A.6.1 Lyapunov exponents are limits of logarithms of subspace spectral norm divided by layer depth $L$	16
A.6.2 Singular value distributions of Jacobian matrices of deep function coupling	17
A.6.3 Existence of Lyapunov exponents for Jacobian matrices of deep function coupling	19
A.7 Proof to Theorem 5	19
A.8 Proof to Corollary 1	19
<b>B Possible Remission Approaches to Rank Diminishing</b>	<b>19</b>
<b>C Code</b>	<b>20</b>
<b>D Partial Rank of Jacobians under Different Input Patches</b>	<b>20</b>
<b>E Estimating Dimension Diminishing in Features</b>	<b>20</b>
<b>F More Examples of Independence Deficit</b>	<b>22</b>

## A Proofs

### A.1 Proof to Lemma 1

This Lemma is the direct result of the *rank theorem of manifolds*.

**Theorem 6** (Rank Theorem [21]). *Suppose  $f : \mathcal{M} \rightarrow \mathcal{N}$  is a smooth function from  $m$ -dimensional manifold  $\mathcal{M}$  to  $n$ -dimensional manifold  $\mathcal{N}$ , and  $\text{Rank}_{\mathcal{M}, \mathcal{N}}(J_f) = r$ . Then for each  $x \in \mathcal{M}$ , there exists a smooth chart  $(U, m)$  around  $x$  and a smooth chart  $(V, n)$  around  $f(x)$ , such that*

$$n \circ f \circ m^{-1} : m(U) \subset \mathbb{R}^m \rightarrow n(V) \subset \mathbb{R}^n \quad (\text{A13})$$

is given by  $n \circ f \circ m^{-1}(x_1, x_2, \dots, x_m) = (x_1, x_2, \dots, x_m, 0, \dots, 0)$ .

The rank of a function  $\text{Rank}_{\mathcal{M}, \mathcal{N}}$  defined on the manifold is the rank under local chart systems of the input manifold  $\mathcal{M}$  and output manifold  $\mathcal{N}$ . Let  $\phi$  be the chart for point  $x \in \mathcal{O} \subset \mathcal{M}$ , and the identity map be the chart for point  $f(x) \in \mathbb{R}^d$ . Then it is easy to find that

$$\text{Rank}(I^{-1} \circ f \circ \phi) = \text{Rank}(f \circ \phi) = \text{Rank}_{\mathcal{M}, \mathcal{N}}(f). \quad (\text{A14})$$

Then by Theorem 6, we know that

$$\dim(f(\mathcal{X})) = \text{Rank}_{\mathcal{M}, \mathcal{N}}(f) = \text{Rank}(f \circ \phi) \leq \text{Rank}(f) = r. \quad (\text{A15})$$

The last equal sign comes from the rank inequality of matrix multiplication  $\text{Rank}(AB) \leq \min\{\text{Rank}(A), \text{Rank}(B)\}$ , which we will discuss later.

### A.2 Proof to Theorem 1

Proof to this Theorem needs Weyl's inequalities [50] for singular values of sum of matrices.

**Theorem 7** (Weyl's inequalities). *Let  $A, B$  be  $p \times n$  complex matrices,  $\sigma_i(\cdot)$  be the  $i$ -th largest singular value of the matrix. Then*

$$|\sigma_i(A + B) - \sigma_i(B)| \leq \sigma_1(B), 1 \leq i \leq p, n. \quad (\text{A16})$$

Let  $p$  be the number of singular values of  $W$  and  $D$ . By this theorem, we have

$$\sigma_i(W) - \delta\sigma_1(D) \leq \sigma_i(W + \delta D) \leq \sigma_i(W) + \delta\sigma_1(D), i = 1, \dots, p. \quad (\text{A17})$$

To measure the numerical rank, we need to estimate the relative quantities of singular values, which are,

$$\frac{\sigma_i(W) - \delta\sigma_1(D)}{\sigma_1(W) + \delta\sigma_1(D)} \leq \frac{\sigma_i(W + \delta D)}{\sigma_1(W + \delta D)} \leq \frac{\sigma_i(W) + \delta\sigma_1(D)}{\sigma_1(W) - \delta\sigma_1(D)}, i = 2, \dots, p. \quad (\text{A18})$$

Now assume that  $\epsilon$  does not belong to the following set (which is a zero measure set in  $\mathbb{R}_+$ )

$$\Sigma_W = \left\{ \frac{\sigma_i(W)}{\sigma_1(W)} : i = 2, \dots, p \right\}, \quad (\text{A19})$$

and  $\text{Rank}_\epsilon(W) = r$ . We know that

$$\frac{\sigma_i(W)}{\sigma_1(W)} > \epsilon, i = 2, \dots, r; \frac{\sigma_i(W)}{\sigma_1(W)} < \epsilon, i = r + 1, \dots, p. \quad (\text{A20})$$

Thus, we have that  $\forall \delta < \delta_{\max}$ ,

$$\frac{\sigma_i(W + \delta D)}{\sigma_1(W + \delta D)} \geq \frac{\sigma_i(W) - \delta\sigma_1(D)}{\sigma_1(W) + \delta\sigma_1(D)} > \epsilon, i = 2, \dots, r, \quad (\text{A21})$$

$$\frac{\sigma_i(W)}{\sigma_1(W)} \leq \frac{\sigma_i(W) + \delta\sigma_1(D)}{\sigma_1(W) - \delta\sigma_1(D)} < \epsilon, i = r + 1, \dots, p, \quad (\text{A22})$$

provided that

$$\begin{aligned} \delta_{\max} = \min\{ & \frac{1}{\sigma_1(D)} \left( \frac{\sigma_r(W) + \sigma_1(W)}{\epsilon + 1} - \sigma_1(W) \right), \frac{\sigma_r(W)}{2\sigma_1(D)}, \\ & \frac{1}{\sigma_1(D)} \left( \sigma_1(W) - \frac{\sigma_{r+1}(W) + \sigma_1(W)}{\epsilon + 1} \right), \frac{\sigma_1(W)}{2\sigma_1(D)} \}. \end{aligned} \quad (\text{A23})$$

Thus we can conclude

$$\text{Rank}_\epsilon(\mathbf{W} + \delta\mathbf{D}) = \text{Rank}_\epsilon(\mathbf{W}), \forall \delta \in [0, \delta_{\max}]. \quad (\text{A24})$$

When  $\text{Rank}(\mathbf{W}) = r < p$ , it is then easy to see if  $\epsilon < \frac{\sigma_r(\mathbf{W})}{\sigma_1(\mathbf{W})}$ , we have

$$\frac{\sigma_i(\mathbf{W})}{\sigma_1(\mathbf{W})} > \epsilon, i = 2, \dots, r; \frac{\sigma_i(\mathbf{W})}{\sigma_1(\mathbf{W})} = 0 < \epsilon, i = r+1, \dots, p. \quad (\text{A25})$$

Thus setting  $\epsilon_{\max} = \frac{\sigma_r(\mathbf{W})}{\sigma_1(\mathbf{W})}$ , we can always have  $\delta_{\max}$  acquired by Eq. (A23), such that

$$\text{Rank}_\epsilon(\mathbf{W}) = \text{Rank}(\mathbf{W}) = \text{Rank}_\epsilon(\mathbf{W} + \delta\mathbf{D}), \forall \delta \in [0, \delta_{\max}]. \quad (\text{A26})$$

### A.3 Proof to Lemma 2

Let  $\mathbf{A}_i$  be the  $i$ -th column of matrix  $\mathbf{A}$ . Given two matrices  $\mathbf{A} \in \mathbb{R}^{m \times n}, \mathbf{B} \in \mathbb{R}^{n \times d}$ , we have

$$\mathbf{AB} = (\mathbf{AB}_1, \dots, \mathbf{AB}_d). \quad (\text{A27})$$

Thus for any  $1 \leq i_1 < \dots < i_K \leq d$ ,

$$(\mathbf{AB})_{i_1, \dots, i_K} = (\mathbf{AB}_{i_1}, \dots, \mathbf{AB}_{i_K}) = \mathbf{A}(\mathbf{B})_{i_1, \dots, i_K}. \quad (\text{A28})$$

By the rank theorem [22] of matrices, we have

$$\text{Rank}((\mathbf{AB})_{i_1, \dots, i_K}) = \text{Rank}((\mathbf{B})_{i_1, \dots, i_K}) - \dim(\text{Ker}(\mathbf{A}) \cap \text{Im}((\mathbf{B})_{i_1, \dots, i_K})), \quad (\text{A29})$$

and

$$\text{Rank}(\mathbf{AB}) = \text{Rank}(\mathbf{B}) - \dim(\text{Ker}(\mathbf{A}) \cap \text{Im}(\mathbf{B})). \quad (\text{A30})$$

As  $(\mathbf{B})_{i_1, \dots, i_K} \subset \mathbf{B}$ , it is straightforward to get that

$$\text{Im}((\mathbf{B})_{i_1, \dots, i_K}) \subset \text{Im}(\mathbf{B}). \quad (\text{A31})$$

Thus

$$\text{Ker}(\mathbf{A}) \cap \text{Im}((\mathbf{B})_{i_1, \dots, i_K}) \subset \text{Ker}(\mathbf{A}) \cap \text{Im}(\mathbf{B}). \quad (\text{A32})$$

Then we have

$$\begin{aligned} \text{Rank}(\mathbf{B}) - \text{Rank}(\mathbf{AB}) &= \dim(\text{Ker}(\mathbf{A}) \cap \text{Im}(\mathbf{B})) \geq \dim(\text{Ker}(\mathbf{A}) \cap \text{Im}((\mathbf{B})_{i_1, \dots, i_K})) \\ &= \text{Rank}((\mathbf{B})_{i_1, \dots, i_K}) - \text{Rank}((\mathbf{AB})_{i_1, \dots, i_K}) \geq 0. \end{aligned} \quad (\text{A33})$$

Note that  $\text{Rank}(\mathbf{f}_2 \circ \mathbf{f}_1) = \text{Rank}(\mathbf{J}_{\mathbf{f}_2} \mathbf{J}_{\mathbf{f}_1})$ . Then we complete the proof.

### A.4 Proof to Theorem 2

The key to this principle is the rank theorem of matrices [22], which is

$$\text{Rank}(\mathbf{AB}) = \text{Rank}(\mathbf{B}) - \dim(\text{Ker}(\mathbf{A}) \cap \text{Im}(\mathbf{B})). \quad (\text{A34})$$

Note that  $\text{Rank}(\mathbf{AB}) = \text{Rank}(\mathbf{B}^T \mathbf{A}^T)$  and  $\text{Rank}(\mathbf{A}^T) = \text{Rank}(\mathbf{A})$ . Then we have

$$\begin{aligned} \text{Rank}(\mathbf{AB}) &= \text{Rank}(\mathbf{B}^T \mathbf{A}^T) = \text{Rank}(\mathbf{A}^T) - \dim(\text{Ker}(\mathbf{B}^T) \cap \text{Im}(\mathbf{A}^T)) \\ &= \text{Rank}(\mathbf{A}) - \dim(\text{Ker}(\mathbf{B}^T) \cap \text{Im}(\mathbf{A}^T)). \end{aligned} \quad (\text{A35})$$

The dimension of a linear subspace will at least be zero, thus the above equations suggest

$$\text{Rank}(\mathbf{AB}) \leq \text{Rank}(\mathbf{B}), \text{Rank}(\mathbf{AB}) \leq \text{Rank}(\mathbf{A}). \quad (\text{A36})$$

Applying this argument to the chain rule of differentials then yields the conclusion. Further using Lemma 1 gives the diminishing of intrinsic dimensions of feature manifolds.

## A.5 Proof to Theorem 3

We first give the rigorous version of this theorem as follows.

**Theorem 8.** Let  $e_x$  be the exponential map from a small neighborhood  $\mathcal{U}_x$  of point  $x$  on the input feature manifold to its tangent space at  $x$ , and  $v_x = e_x(x)$ . Let  $\mathcal{X}$  be the input manifold,  $s = \dim(\mathcal{X})$ ,  $f^i$  be the layer network,  $r = \text{Rank}(f^i)$ , and  $f^i(\mathcal{X})$  be the output manifold. If for almost everywhere on the input feature manifold, there is a unit vector  $v \in e_x(\mathcal{U}_x)$ , such that the layer network  $f^i$  satisfies

$$\lim_{t \rightarrow 0} \frac{\|f^i \circ e_x^{-1}(v_x) - f^i \circ e_x^{-1}(v_x + tv)\|_2}{t} = 0, \quad (\text{A37})$$

then  $\dim(f^i(\mathcal{X})) < s$ . If the number of such independent  $v$  in  $e_x(\mathcal{U}_x)$  is  $k$ , then  $\dim(f^i(\mathcal{X})) \leq s - k$ .

Now we prove this theorem.

Note that Eq. (A37) implies

$$J_{e_x^{-1}} v \in \text{Ker}(J_{f^i}). \quad (\text{A38})$$

As it is also easy to see

$$J_{e_x^{-1}} v \in \text{Im}(J_{e_x^{-1}}), \quad (\text{A39})$$

we can conclude

$$\mathbf{0} \neq J_{e_x^{-1}} v \in \text{Ker}(J_{f^i}) \cap \text{Im}(J_{e_x^{-1}}), \quad (\text{A40})$$

where  $\mathbf{0} \neq J_{e_x^{-1}} v$  comes from the full rank property of exponential map and its inverse. Thus we have

$$\dim(\text{Ker}(J_{f^i}) \cap \text{Im}(J_{e_x^{-1}})) \geq 1. \quad (\text{A41})$$

Specifically, if linearly independent  $v_1, \dots, v_k$  satisfy Eq. (A37), we can conclude

$$\mathbf{0} \neq J_{e_x^{-1}} v_i \in \text{Ker}(J_{f^i}) \cap \text{Im}(J_{e_x^{-1}}), i = 1, \dots, k. \quad (\text{A42})$$

As  $J_{e_x^{-1}}$  has full rank due to the property of exponential map, we know that  $J_{e_x^{-1}} v_i, i = 1, \dots, k$  are linearly independent. Then

$$\dim(\text{Ker}(J_{f^i}) \cap \text{Im}(J_{e_x^{-1}})) \geq k. \quad (\text{A43})$$

Thus the rank theorem of matrices [22] reads

$$\text{Rank}(J_{f^i \circ e_x^{-1}}) = \text{Rank}(J_{f^i} J_{e_x^{-1}}) = \text{Rank}(J_{e_x^{-1}}) - \dim(\text{Ker}(J_{f^i}) \cap \text{Im}(J_{e_x^{-1}})) \quad (\text{A44})$$

$$= s - \dim(\text{Ker}(J_{f^i}) \cap \text{Im}(J_{e_x^{-1}})) \leq s - k. \quad (\text{A45})$$

Combining this result with Theorem 6 proves our result.

## A.6 Proof to Theorem 4

The proof to this theorem relies on the existence of Lyapunov exponents of dynamic systems. Given a linearized dynamic system

$$\dot{v}(t) = X_t v, v(0) = v_0 \in \mathbb{R}^n, \quad (\text{A46})$$

its (largest) Lyapunov exponent is defined as

$$\lambda = \limsup_{t \rightarrow \infty} \frac{1}{t} \|v\|_2. \quad (\text{A47})$$

Further, for a sequence of subspace  $\mathcal{L}_h \subset \mathcal{L}_{r-1} \subset \dots \subset \mathcal{L}_1 \subset \mathcal{L}_0 = \mathbb{R}^n$ , we can define the corresponding Lyapunov exponents of all those subspaces as

$$\lambda_i = \lim_{t \rightarrow \infty} \frac{1}{t} \log \|v\|_2, i = 1, \dots, h+1, v_0 \in \mathcal{L}_{i-1} \setminus \mathcal{L}_i, \quad (\text{A48})$$

and we have

$$\lambda = \lambda_1 > \lambda_2 > \dots > \lambda_h. \quad (\text{A49})$$

It may be surprising to find that such Lyapunov exponents exist, as  $\mathbf{v}_0$  can traverse the entire subspace  $\mathcal{L}_{i-1} \setminus \mathcal{L}_i$ . We will demonstrate the existence of the Lyapunov exponents for our case later in Sec. A.6.3, which is the classical results from the Furstenberg-Kesten theorem [16] and multiplicative ergodic theorem [41]. Before that, we will first assume the existence of those Lyapunov exponents for simplicity of analysis.

Now consider the case of function couplings

$$\mathbf{F} = \mathbf{f}^L \circ \cdots \circ \mathbf{f}^2 \circ \mathbf{f}^1, \quad (\text{A50})$$

which has the Jacobian matrix

$$\mathbf{J}_F = \mathbf{J}_{\mathbf{f}^L} \mathbf{J}_{\mathbf{f}^{L-1}} \cdots \mathbf{J}_{\mathbf{f}^2} \mathbf{J}_{\mathbf{f}^1}. \quad (\text{A51})$$

Apparently, the following dynamic system induces the Jacobi matrix of  $\mathbf{F}$ ,

$$\dot{\mathbf{v}}(t) = \mathbf{J}_{\mathbf{f}^t} \mathbf{v}, \quad \mathbf{v}(0) = \mathbf{v}_0 \in \mathbb{R}^n, \quad t = 1, \dots, L. \quad (\text{A52})$$

Thus its Lyapunov exponents are given by

$$\lambda_i = \lim_{L \rightarrow \infty} \frac{1}{L} \log \|\mathbf{J}_F \mathbf{v}_0\|_2, \quad i = 1, \dots, h+1, \quad \mathbf{v}_0 \in \mathcal{L}_{i-1} \setminus \mathcal{L}_i, \quad (\text{A53})$$

for a chain of subspaces  $\{0\} = \mathcal{L}_{h+1} \subset \mathcal{L}_h \subset \mathcal{L}_{h-1} \subset \cdots \subset \mathcal{L}_1 \subset \mathcal{L}_0 = \mathbb{R}^n$ .

### A.6.1 Lyapunov exponents are limits of logarithms of subspace spectral norm divided by layer depth $L$

We first demonstrate that the Lyapunov exponents are limits of logarithm of the spectral norm of  $\mathbf{F}$  on  $\mathcal{L}_{i-1} \setminus \mathcal{L}_i$  divided by layer depth  $L$  when  $L \rightarrow \infty$ , for  $i = 1, \dots, r$ .

It is easy to see

$$\frac{1}{L} \log \|\mathbf{v}_0\|_2 \|\mathbf{J}_F \frac{\mathbf{v}_0}{\|\mathbf{v}_0\|_2}\|_2 = \frac{1}{L} (\log \|\mathbf{v}_0\|_2 + \log \|\mathbf{J}_F \frac{\mathbf{v}_0}{\|\mathbf{v}_0\|_2}\|_2). \quad (\text{A54})$$

When  $L \rightarrow \infty$ ,  $\frac{1}{L} \log \|\mathbf{v}_0\|_2 \rightarrow 0$  for any  $\mathbf{v}_0$ , we have

$$\lambda_i = \lim_{L \rightarrow \infty} \frac{1}{L} \log \|\mathbf{J}_F \mathbf{v}_0\|_2, \quad \|\mathbf{v}_0\|_2 = 1, \quad \mathbf{v}_0 \in \mathcal{L}_{i-1} \setminus \mathcal{L}_i. \quad (\text{A55})$$

Let

$$\lambda_i^L = \sup_{\|\mathbf{v}_0\|_2=1, \mathbf{v}_0 \in \mathcal{L}_{i-1} \setminus \mathcal{L}_i} \frac{1}{L} \log \|\mathbf{J}_F \mathbf{v}_0\|_2. \quad (\text{A56})$$

Note that

$$\sup_{\|\mathbf{v}_0\|_2=1, \mathbf{v}_0 \in \mathcal{L}_{i-1} \setminus \mathcal{L}_i} \frac{1}{L} \log \|\mathbf{J}_F \mathbf{v}_0\|_2 = \frac{1}{L} \log \sup_{\|\mathbf{v}_0\|_2=1, \mathbf{v}_0 \in \mathcal{L}_{i-1} \setminus \mathcal{L}_i} \|\mathbf{J}_F \mathbf{v}_0\|_2, \quad (\text{A57})$$

and

$$\sup_{\|\mathbf{v}_0\|_2=1, \mathbf{v}_0 \in \mathcal{L}_{i-1} \setminus \mathcal{L}_i} \|\mathbf{J}_F \mathbf{v}_0\|_2 = \|\mathbf{J}_F\|_{2,i}, \quad (\text{A58})$$

where  $\|\cdot\|_{2,i}$  denote the spectral norm of a linear operator constrained on  $\mathcal{L}_{i-1} \setminus \mathcal{L}_i$ . Then we have

$$\lambda_i^L = \frac{1}{L} \log \|\mathbf{J}_F\|_{2,i}. \quad (\text{A59})$$

Let  $\mathbf{e}_{i_1}, \dots, \mathbf{e}_{i_k}$  be a set of standard orthogonal basis of  $i_k$  dimensional subspace  $\mathcal{L}_{i-1} \setminus \mathcal{L}_i$ . If the Lyapunov exponents exist, by Eq. (A55) we have for any  $\epsilon > 0$ , there is  $N \in \mathbb{N}$  such that for all  $L > N$ ,

$$\lambda_i - \frac{\epsilon}{2} \leq \frac{1}{L} \log \|\mathbf{J}_F \mathbf{e}_j\|_2 \leq \lambda_i + \frac{\epsilon}{2}, \quad j = 1, \dots, i_k. \quad (\text{A60})$$

Let  $\mathbf{v} = \sum_{j=1}^{i_k} \alpha_j \mathbf{e}_j \in \mathcal{L}_{i-1} \setminus \mathcal{L}_i$ , where  $\alpha_j, j = 1, \dots, i_k, \sum_{j=1}^{i_k} \alpha_j^2 = 1$  is the coordinate of unit vector  $\mathbf{v}$  under the basis  $\mathbf{e}_j, j = 1, \dots, i_k$ . Assume that  $\|\mathbf{J}_F \mathbf{e}_1\|_2 \geq \|\mathbf{J}_F \mathbf{e}_j\|_2, j = 2, \dots, i_k$ . We then have  $\frac{|\alpha_j| \|\mathbf{J}_F \mathbf{e}_j\|_2}{\sum_{j=1}^{i_k} |\alpha_j| \|\mathbf{J}_F \mathbf{e}_1\|_2} \leq 1, j = 1, \dots, i_k$ , and

$$\frac{1}{L} \log \|\mathbf{J}_F \mathbf{v}\|_2 \leq \frac{1}{L} \log \sum_{j=1}^{i_k} |\alpha_j| \|\mathbf{J}_F \mathbf{e}_j\|_2 \quad (\text{A61})$$

$$= \frac{1}{L} \log \left( \frac{|\alpha_1| \|\mathbf{J}_F \mathbf{e}_1\|_2}{\sum_{j=1}^{i_k} |\alpha_j| \|\mathbf{J}_F \mathbf{e}_1\|_2} + \sum_{j=2}^{i_k} \frac{|\alpha_j| \|\mathbf{J}_F \mathbf{e}_j\|_2}{\sum_{j=1}^{i_k} |\alpha_j| \|\mathbf{J}_F \mathbf{e}_1\|_2} \right) \left( \sum_{j=1}^{i_k} |\alpha_j| \|\mathbf{J}_F \mathbf{e}_1\|_2 \right) \quad (\text{A62})$$

$$= \frac{1}{L} \log \left( \frac{|\alpha_1| \|\mathbf{J}_F \mathbf{e}_1\|_2}{\sum_{j=1}^{i_k} |\alpha_j| \|\mathbf{J}_F \mathbf{e}_1\|_2} + \sum_{j=2}^{i_k} \frac{|\alpha_j| \|\mathbf{J}_F \mathbf{e}_j\|_2}{\sum_{j=1}^{i_k} |\alpha_j| \|\mathbf{J}_F \mathbf{e}_1\|_2} \right) + \frac{1}{L} \log \sum_{j=1}^{i_k} |\alpha_j| \|\mathbf{J}_F \mathbf{e}_1\|_2 \quad (\text{A63})$$

$$\leq \frac{1}{L} \log \left( 1 + \sum_{j=2}^{i_k} \frac{|\alpha_j| \|\mathbf{J}_F \mathbf{e}_j\|_2}{\sum_{j=1}^{i_k} |\alpha_j| \|\mathbf{J}_F \mathbf{e}_1\|_2} \right) + \frac{1}{L} \log \sum_{j=1}^{i_k} |\alpha_j| \|\mathbf{J}_F \mathbf{e}_1\|_2 \quad (\text{A64})$$

$$\leq \frac{1}{L} \sum_{j=2}^{i_k} \frac{|\alpha_j| \|\mathbf{J}_F \mathbf{e}_j\|_2}{\sum_{j=1}^{i_k} |\alpha_j| \|\mathbf{J}_F \mathbf{e}_1\|_2} + \frac{1}{L} \log \sum_{j=1}^{i_k} |\alpha_j| \|\mathbf{J}_F \mathbf{e}_1\|_2 \quad (\text{A65})$$

$$\leq \frac{1}{L} (i_k - 1) + \frac{1}{L} \log \sum_{j=1}^{i_k} |\alpha_j| + \frac{1}{L} \log \|\mathbf{J}_F \mathbf{e}_1\|_2 \quad (\text{A66})$$

$$\leq \frac{1}{L} (i_k - 1) + \frac{1}{2L} \log(1^2 + \dots + 1^2)(\alpha_1^2 + \dots + \alpha_{i_k}^2) + \frac{1}{L} \log \|\mathbf{J}_F \mathbf{e}_1\|_2 \quad (\text{A67})$$

$$= \frac{1}{L} (i_k - 1) + \frac{1}{2L} \log i_k + \frac{1}{L} \log \|\mathbf{J}_F \mathbf{e}_1\|_2 \quad (\text{A68})$$

$$\leq \frac{1}{L} (i_k - 1) + \frac{1}{2L} \log i_k + \lambda_i + \frac{\epsilon}{2}. \quad (\text{A69})$$

Thus, if we set  $N_0 = \max\{N, \frac{2i_k - 2 + \log i_k}{\epsilon}\}$ , then when  $L > N_0$  we have  $\frac{1}{L} (i_k - 1) + \frac{1}{2L} \log i_k < \frac{\epsilon}{2}$  and

$$\frac{1}{L} \log \|\mathbf{J}_F \mathbf{v}\|_2 \leq \lambda_i + \epsilon, \quad \forall \mathbf{v} \in \mathcal{L}_{i-1} \setminus \mathcal{L}_i, \quad \|\mathbf{v}\|_2 = 1. \quad (\text{A70})$$

Combining Eqs. (A60) and (A70), we have for any  $\epsilon > 0$ , there is  $N_0 \in \mathbb{N}$ , such that when  $L > N_0$ , we always have

$$\lambda_i - \frac{\epsilon}{2} \leq \lambda_i^L = \sup_{\|\mathbf{v}_0\|_2=1, \mathbf{v}_0 \in \mathcal{L}_{i-1} \setminus \mathcal{L}_i} \frac{1}{L} \log \|\mathbf{J}_F \mathbf{v}_0\|_2 = \frac{1}{L} \log \|\mathbf{J}_F\|_{2,i} \leq \lambda_i + \epsilon. \quad (\text{A71})$$

Thus, if the Lyapunov exponents exist, *i.e.*, the existence of limits of Eq. (A53), we have

$$\lambda_i = \lim_{L \rightarrow \infty} \frac{1}{L} \log \|\mathbf{J}_F\|_{2,i} = \lim_{L \rightarrow \infty} \lambda_i^L. \quad (\text{A72})$$

### A.6.2 Singular value distributions of Jacobian matrices of deep function coupling

In Sec. A.6.1 we have proved that the Lyapunov exponents (if they exist) are limits of logarithms of subspace spectral norms divided by  $L$ . Here we use this property to prove the deficiency of numerical ranks, *i.e.*, Eq. (9).

We first introduce the Courant-Fischer min-max theorem [22] of singular values.

**Theorem 9** (Courant-Fischer Min-max Theorem). *Let  $\mathbf{A}$  be a  $d \times n$  complex matrix and  $\sigma_i(\mathbf{A})$  denote its  $i$ -th largest singular value,  $i = 1, \dots, \min\{d, n\}$ . Then we have*

$$\sigma_i(\mathbf{A}) = \sup_{\dim(\mathcal{V})=i} \inf_{\mathbf{v} \in \mathcal{V}, \|\mathbf{v}\|_2=1} \|\mathbf{A}\mathbf{v}\|_2, \quad (\text{A73})$$

$$\sigma_i(\mathbf{A}) = \inf_{\dim(\mathcal{V})=n-i+1} \sup_{\mathbf{v} \in \mathcal{V}, \|\mathbf{v}\|_2=1} \|\mathbf{A}\mathbf{v}\|_2, \quad (\text{A74})$$

$$(A75)$$

where  $\mathcal{V}$  traverses subspaces of  $\mathbb{R}^n$ .

This theorem also serves as one of the definitions to singular values.

Now assume that  $\dim(\mathcal{L}_0 \setminus \mathcal{L}_1) = r$ , and we consider only the case of  $d = n$  for simplicity. For any  $\epsilon > 0$ , we have  $N \in \mathbb{N}$  such that when  $L > N$ ,

$$\inf_{\mathbf{v} \in \mathcal{L}_0 \setminus \mathcal{L}_1, \|\mathbf{v}\|_2=1} \|\mathbf{J}_F \mathbf{v}\|_2 \geq \exp L(\lambda_1 - \epsilon) \quad (\text{A76})$$

$$(\text{A77})$$

due to Eq. (A55). As  $\dim(\mathcal{L}_0 \setminus \mathcal{L}_1) = r$ , by Theorem 9, we have

$$\sigma_1(\mathbf{J}_F) \geq \dots \geq \sigma_r(\mathbf{J}_F) \geq \exp L(\lambda_1 - \epsilon). \quad (\text{A78})$$

For  $\mathbf{v} \in \mathcal{L}_0 = \mathbb{R}^n$  and  $\|\mathbf{v}\|_2 = 1$ , as  $\mathcal{L}_0 = \mathcal{L}_0 \setminus \mathcal{L}_1 \oplus \dots \oplus \mathcal{L}_h \setminus \mathcal{L}_{h+1}$  ( $\oplus$  denotes direct sum of linear quotient subspaces in the Banach space), there is  $\mathbf{v}_i \in \mathcal{L}_{i-1} \setminus \mathcal{L}_i$ ,  $i = 1, \dots, h+1$ , such that

$$\|\mathbf{v}_1\|_2^2 + \dots + \|\mathbf{v}_{h+1}\|_2^2 = 1 \quad (\text{A79})$$

and

$$\mathbf{v} = \mathbf{v}_1 + \dots + \mathbf{v}_{h+1}. \quad (\text{A80})$$

Then by the conclusion of Sec. A.6.1, we have

$$\limsup_{L \rightarrow \infty} \|\mathbf{J}_F \mathbf{v}\|_2 \leq \limsup_{L \rightarrow \infty} \|\mathbf{J}_F \mathbf{v}_1\|_2 + \dots + \limsup_{L \rightarrow \infty} \|\mathbf{J}_F \mathbf{v}_{h+1}\|_2 \quad (\text{A81})$$

$$\leq \|\mathbf{v}_1\|_2 \lim_{L \rightarrow \infty} \|\mathbf{J}_F\|_{2,2} + \dots + \|\mathbf{v}_{h+1}\|_2 \lim_{L \rightarrow \infty} \|\mathbf{J}_F\|_{2,h+1} \leq \sum_{i=1}^{h+1} \|\mathbf{v}_i\|_2 \lambda_i \leq \lambda_1. \quad (\text{A82})$$

Thus there is  $N_1 \in \mathbb{N}$ , such that when  $L > N_1$ , we have

$$\sup_{\mathbf{v} \in \mathcal{L}_0, \|\mathbf{v}\|_2=1} \|\mathbf{J}_F \mathbf{v}\|_2 \leq \lambda_1 + \epsilon. \quad (\text{A83})$$

As  $\dim(\mathcal{L}_1) = n - 1 + 1$ , by Theorem 9, we have when  $L > N_1$ ,

$$\sigma_r(\mathbf{J}_F) \leq \dots \leq \sigma_1(\mathbf{J}_F) \leq \exp L(\lambda_1 + \epsilon). \quad (\text{A84})$$

In conclusion, when  $L > N_0 = \max\{N, N_1\}$ , we have

$$\exp L(\lambda_1 - \epsilon) \leq \sigma_1(\mathbf{J}_F) \leq \exp L(\lambda_1 + \epsilon). \quad (\text{A85})$$

Thus when  $L \rightarrow \infty$ , we have

$$\sigma_1(\mathbf{J}_F) \sim \exp L \lambda_1. \quad (\text{A86})$$

Using the same argument for  $\sigma_2(\mathbf{J}_F), \dots, \sigma_n(\mathbf{J}_F)$ , we can find that if let  $\hat{\lambda}_1 \geq \hat{\lambda}_2 \dots \geq \hat{\lambda}_n$  be the Lyapunov exponents counting repetitions, i.e.,

$$\hat{\lambda}_k = \lambda_i, \text{ if } \sum_{j=1}^{i-1} \dim(\mathcal{L}_{j-1} \setminus \mathcal{L}_j) < k \leq \sum_{j=1}^i \dim(\mathcal{L}_{j-1} \setminus \mathcal{L}_j), i = 1, \dots, h+1, \quad (\text{A87})$$

then

$$\sigma_i(\mathbf{J}_F) \sim \exp L \hat{\lambda}_i. \quad (\text{A88})$$

Note that

$$\hat{\lambda}_1 = \dots = \hat{\lambda}_r = \lambda_1, \hat{\lambda}_i \leq \lambda_2 < \lambda_1, i = r+1, \dots, n. \quad (\text{A89})$$

Thus we have

$$\frac{\sigma_i(\mathbf{J}_F)}{\sigma_1(\mathbf{J}_F)} \sim \exp L(\hat{\lambda}_i - \hat{\lambda}_1) \rightarrow 0, i = r+1, \dots, n. \quad (\text{A90})$$

As a consequence,  $\text{Rank}_\epsilon(\mathbf{F}) \leq r$  for any  $\epsilon > 0$  when  $L \rightarrow \infty$ .

### A.6.3 Existence of Lyapunov exponents for Jacobian matrices of deep function coupling

In above analysis, we have proven Theorem 4 under the existence of Lyapunov exponents. In this section, we introduce the classical result of multiplicative ergodic theorem in the specific domain of random matrices, which is proposed by Furstenberg and Kesten [16, 17].

**Theorem 10** (Multiplicative Ergodic Theorem (Theorem 3.9 of [17])). *Let  $\mu$  be a probability measure on all convertible matrices of  $\mathbb{R}^{n \times n}$  which satisfies*

$$\mathbb{E}_\mu[\max\{\log \|J_{f^k}^{\pm 1}\|_2, 0\}] < \infty, k = 1, \dots, L. \quad (\text{A91})$$

*If each  $J_{f^k}$  independently follows  $\mu$ , then we have a chain of subspaces  $\{0\} = \mathcal{L}_{h+1} \subset \mathcal{L}_h \subset \dots \subset \mathcal{L}_1 \subset \mathcal{L}_0 = \mathbb{R}^n$  and corresponding positive real constants  $\lambda_1 > \lambda_2 > \dots > \lambda_{h+1}$  such that almost surely*

$$\lambda_i = \lim_{t \rightarrow 0} \frac{1}{t} \log \|J_F v\|_2, \forall v \in \mathcal{L}_{i-1} \setminus \mathcal{L}_i, i = 1, \dots, h+1, \quad (\text{A92})$$

*which means the existence of the Lyapunov exponents.*

Combining this theorem and the arguments above, we can finally prove Theorem 4.

### A.7 Proof to Theorem 5

This theorem can be deduced from the Lyapunov components of Ginibre matrices (polynomial ensemble of square matrices sampled *i.i.d* from standard Gaussian).

**Theorem 11** (Exact Lyapunov Exponent Distribution for Ginibre Matrices [38]). *If  $\mu$  in Theorem 10 is standard Gaussian, then  $h+1 = n$ , and*

$$\lambda_i = \log \left( 2 + \psi\left(\frac{n-i+1}{2}\right) \right), i = 1, \dots, n. \quad (\text{A93})$$

Combining this theorem with Theorem 4 can directly yield our result.

### A.8 Proof to Corollary 1

This theorem is the direct result of Theorems 4 and 11. Note that it is easy to get

$$\lambda_1 = \lim_{L \rightarrow \infty} \frac{1}{L} \log \|J_F\|_2 \quad (\text{A94})$$

for standard Gaussian  $\mu$ .

## B Possible Remission Approaches to Rank Diminishing

**Skip Connection** Skip Connection is the most direct method to solve rank diminishing. In our formulation, the definition of a layer network requires it to accept inputs purely from its predecessor layer as

$$x^i = f^i(x^{i-1}), x^{i-1} = f^{i-1}(x^{i-2}). \quad (\text{A95})$$

However, when we add a skip connection from its ancestor layer  $f^s, s < i-1$ , we have

$$x^i = f^i(x^{i-1}, x^s), x^{i-1} = f^{i-1}(x^{i-2}), x^{i-2} = f^{i-2} \circ \dots \circ f^s(x^{s-1}), x^s = f^s(x^{s-1}). \quad (\text{A96})$$

It actually makes the coupling of layers

$$\hat{f}^s = \begin{pmatrix} f^{i-1} \circ \dots \circ f^s \\ f^s \end{pmatrix}, \quad (\text{A97})$$

the true predecessor layer to  $f^i$ , as

$$x^i = f^i(\hat{x}), \hat{x} = \hat{f}^s(x^{s-1}). \quad (\text{A98})$$

Thus the true layer depth is cut down by  $i-s$ , remaining  $L-(i-s)$  layers. Skip connection is usually used with the residual network. This structure can ease rank diminishing inside the layer  $\hat{f}^s$ , which we will discuss later. Overall, skip connection shortens the length of the chain of Jacobian matrices, thus restraining rank diminishing.

**BatchNorm** Some previous works [12, 4] discuss the role of BatchNorm in restraining rank diminishing. They show that BatchNorm may slow down the speed of rank diminishing in neural networks in some specific cases.

**Residual Network** Residual Network is another useful tool to restrain rank diminishing. The residual network  $r$  has the form

$$\mathbf{x}^o = \mathbf{r}(\mathbf{x}^i) = \mathbf{x}^i + \text{Res}(\mathbf{x}^i), \quad (\text{A99})$$

where  $\mathbf{x}^o$  and  $\mathbf{x}^i$  are the output feature and input feature, respectively. Usually the residual term  $\text{Res}(\mathbf{x}^i)$  is small compared with the input  $\mathbf{x}^i$ . Assume that

$$\|\mathbf{J}_{\text{Res}}\|_2 < \epsilon, \quad (\text{A100})$$

where  $\epsilon$  is very small. Then we have

$$\mathbf{J}_r = \mathbf{I} + \mathbf{J}_{\text{Res}} \quad (\text{A101})$$

is a diagonally dominant matrix, thus it has full rank. This means its kernel space  $\text{Ker}(\mathbf{J}_r) = \{0\}$  is a zero dimension space. Thus by the Rank Theorem, for any predecessor layer  $f$ ,  $r \circ f$  will not lose rank as

$$\begin{aligned} \text{Rank}(r \circ f) &= \text{Rank}(\mathbf{J}_r \mathbf{J}_f) = \text{Rank}(\mathbf{J}_f) - \dim(\text{Ker}(\mathbf{J}_r) \cap \text{Im}(\mathbf{J}_f)) \\ &= \text{Rank}(\mathbf{J}_f) = \text{Rank}(f). \end{aligned} \quad (\text{A102})$$

## C Code

Algorithm A1 provides the pseudo-code of partial rank of the Jacobian. The implementation of the Algorithm A1 can refer to the ‘rank\_jacobian.py’ python file.

Algorithm A2 provides the pseudo-code of perturbed PCA dimension of feature spaces. The implementation of the Algorithm A2 can refer to the ‘rank\_perturb.py’ python file.

Algorithm A3 provides the pseudo-code of the classification dimension. The implementation of the Algorithm A3 can refer to the ‘run\_cls\_dim.py’ python file.

Algorithm A4 provides the pseudo-code of independence deficit. The implementation of the Algorithm A4 can refer to the ‘run\_deficit.py’ python file.

## D Partial Rank of Jacobians under Different Input Patches

In Fig. A5 we report partial ranks of different input image patches (marked with colored boxes in Fig. A5(a)) for the layers of ResNet-50 on ImageNet. We can find that the curves of partial ranks share a similar and consistent trend among different input patches. Thus, picking one patch, for example, the central patch of  $16 \times 16 \times 3$  pixels we use in Sec. 4.3, could be enough to demonstrate the overall behavior of network ranks. The consistent behavior of all those partial ranks also shows that partial rank is a good tool to investigate network ranks.

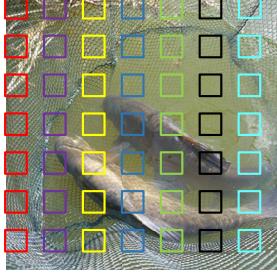
## E Estimating Dimension Diminishing in Features

Measuring the intrinsic dimension of feature manifolds is known to be hard. However, we manage to give a rough estimation to the dimension dropped by different layer networks. To do this, we use a new metric called the Perturbed PCA Dimension. It measures the expectation of PCA dimension of small local neighborhoods over the feature manifold.

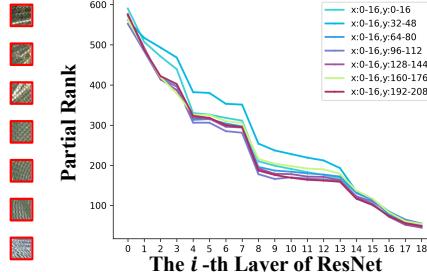
Let  $F_k$  be the  $k$ -th sub-network of the whole network  $F$ . We want to measure the Perturbed PCA Dimension of  $F_k(\mathcal{X})$ , where  $\mathcal{X}$  is the input data domain. To this end, we compute

$$\text{PertDim} = \mathbb{E}_{\mathbf{x} \sim \mathcal{P}_{\mathcal{X}}} [\text{PCADim}(\{F_k(\mathbf{x} + \boldsymbol{\epsilon}) : \boldsymbol{\epsilon} \sim \mathcal{N}(\mathcal{O}, \delta \mathcal{I})\})], \quad (\text{A103})$$

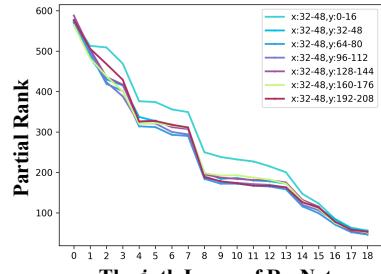
where PCADim for a set is the number of PCA eigenvalues larger than a threshold  $\xi$ . For each point  $\mathbf{x}$ , we sample 50,000 different perturbation  $\boldsymbol{\epsilon}$  to compute the PCA dimension of the neighborhood of  $F_k(\mathbf{x})$ . When computing the PCA dimension, we set  $\delta = 1e-3$  and  $\xi = 1.19e-7 \times 50000 \times \text{eig}_{\max}$ ,



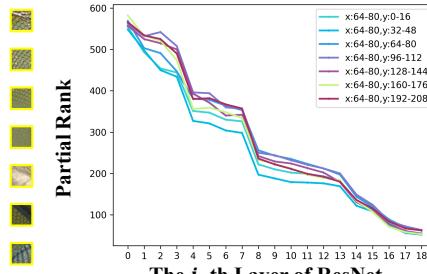
(a) Image



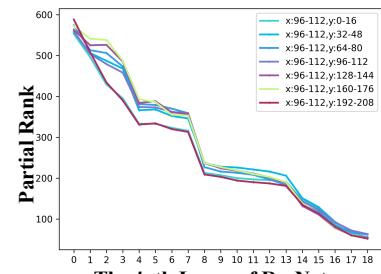
(b) Partial rank in the first column of the image



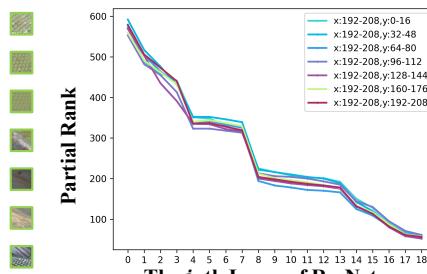
(c) Partial rank in the third column of the image



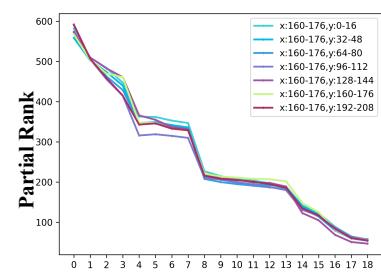
(d) Partial rank in the fifth column of the image



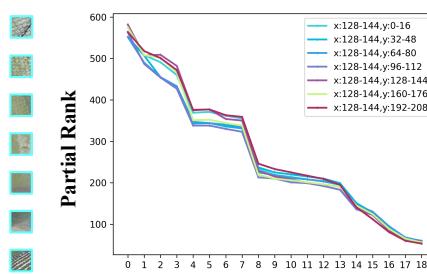
(e) Partial rank in the seventh column of the image



(f) Partial rank in the ninth column of the image



(g) Partial rank in the eleventh column of the image



(h) Partial rank in the thirteenth column of the image

Figure A5: The partial ranks of different input patches at the  $i$ -th layer of ResNet-50 on ImageNet.

where  $\text{eig}_{\max}$  is the largest PCA eigenvalue. We then compute the mean value of PCA dimensions over the neighborhood of 100 random samples in the validation set of ImageNet as the final result.

We do not use PCA dimension of the feature manifolds directly as it is unable to cope with the highly non-linear structure of intermediate feature manifolds. However, the Perturbed PCA Dimension is able to estimate the dimensions of local neighborhoods of points in the feature manifolds. As local neighborhoods can be viewed as linear if the network is smooth, the Perturbed PCA Dimension could be more feasible than PCA dimension in our case. We provide the pseudo-code to compute the Perturbed PCA Dimension in Algorithm A2.

However, the perturbation is made in the ambient space of the input data manifold  $\mathcal{X}$  rather than the data manifold itself. Thus this estimation may considerably overestimate the intrinsic dimensions of feature manifolds. So we merely care about how many Perturbed PCA Dimensions are lost by a sub-network instead of its own Perturbed PCA Dimension. We call this quantity  $\Delta$  Dimension, which is the difference between the Perturbed PCA Dimension of the current layer and that of the input layer for the given deep network. As shown in Fig. A6, we show the dropped dimensions of different feature layers of the CNN, MLP, and Transformer architectures on ImageNet. The results show that the Perturbed PCA Dimensions of feature manifolds of most networks decrease as the networks get deeper, thus confirming the rank diminishing principle we propose in Theorem 2.

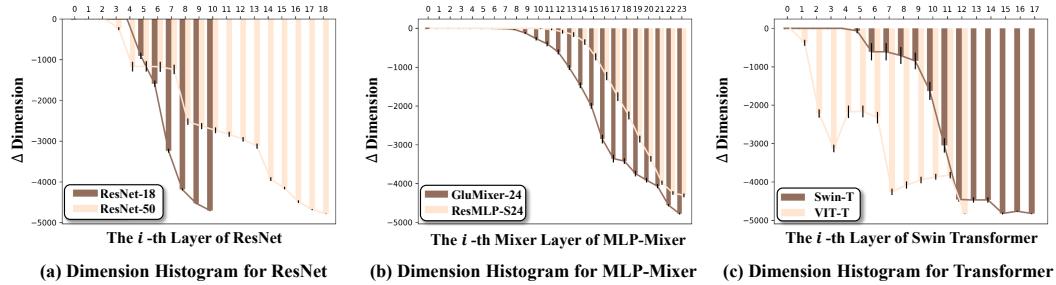


Figure A6: Dropped Perturbed PCA Dimension of different layers.  $\Delta$  Dimension for the  $i$ -th layer is the difference between the Perturbed PCA Dimension of the  $i$ -th layer and that of the input layer of the CNN, MLP, and Transformer architectures on ImageNet.

## F More Examples of Independence Deficit

We provide more examples of independence deficit, which shows that classification confidence of some categories can be linearly decided by a few other categories with fixed coefficients. The results obtained by ResNet-18, ResNet-50, GluMixer-24, ResMLP-S24, ViT-T, and Swin-T are reported in Fig. A7 - Fig. A12, respectively. All the results are obtained by solving the Lasso problem in Sec. 5. In each figure, we report the classification accuracy for category ‘ $i$ ’: the accuracy by calculating logits with Eq. (12) is reported as ‘acc.’; and the original model accuracy is reported as ‘ori. acc.’. Both the metrics are measured in the whole ImageNet validation set. We further report the classification accuracy on positive samples only for both metrics as ‘pos’ following ‘acc.’ and ‘ori. acc.’ correspondingly. The results show a universal independence deficit phenomenon for broad categories in all those deep networks.

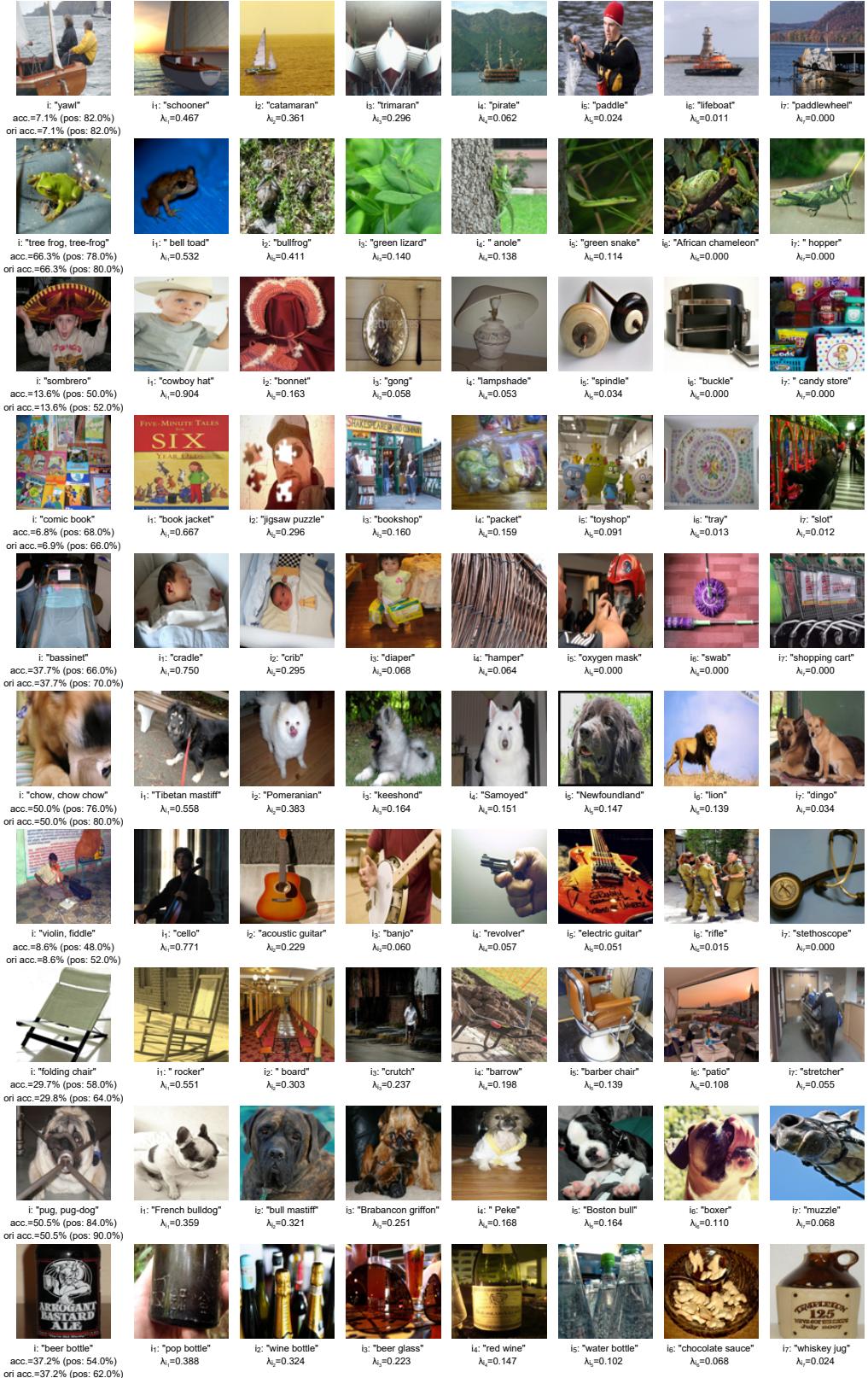


Figure A7: Results from ResNet-18, where ‘acc.’ and ‘ori acc.’ denote the classification accuracies on the ImageNet validation set, while ‘pos: xx%’ is the accuracy on positive samples only.

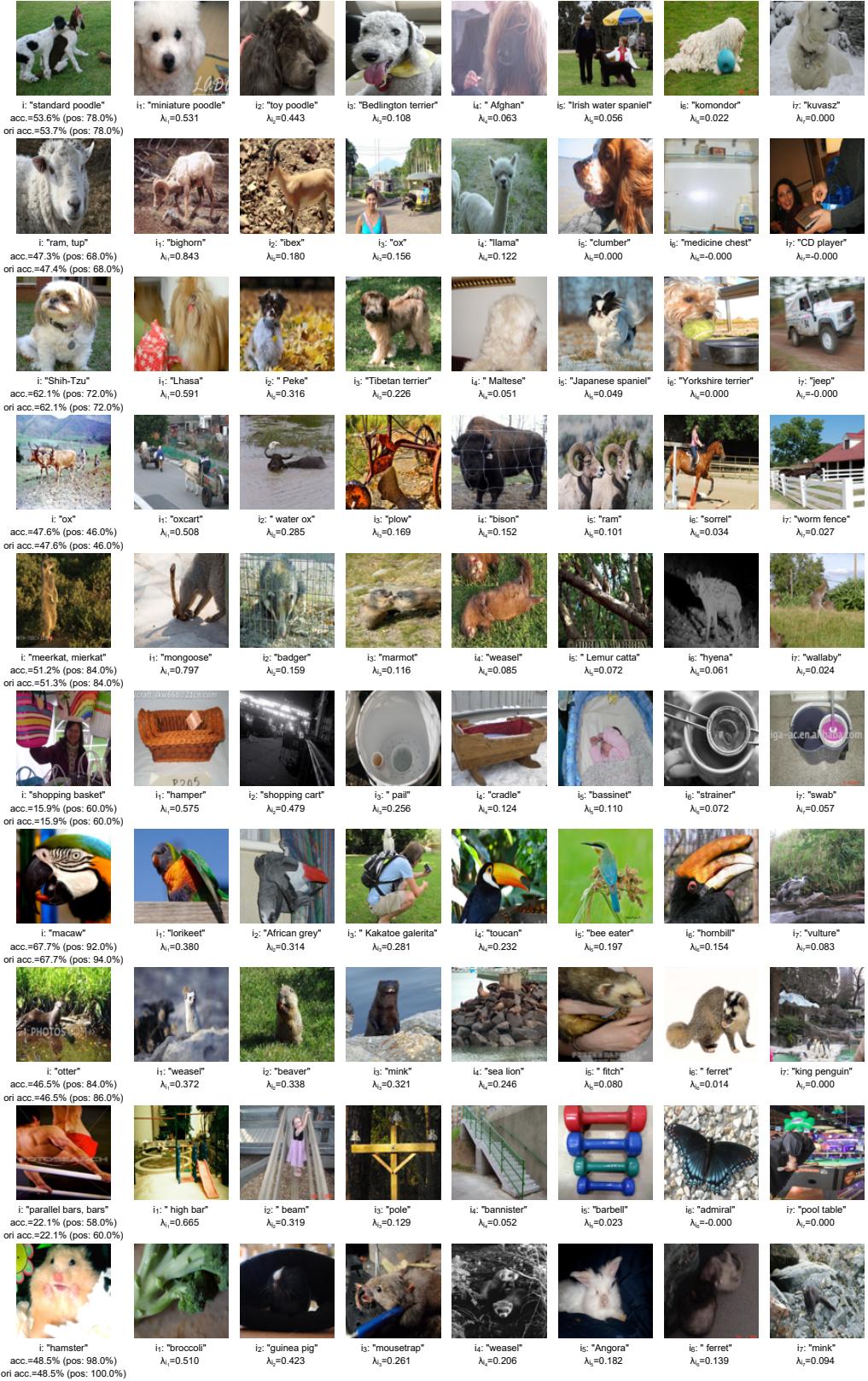


Figure A8: Results from ResNet-50, where ‘acc.’ and ‘ori acc.’ denote the classification accuracies on the ImageNet validation set, while ‘pos: xx%’ is the accuracy on positive samples only.

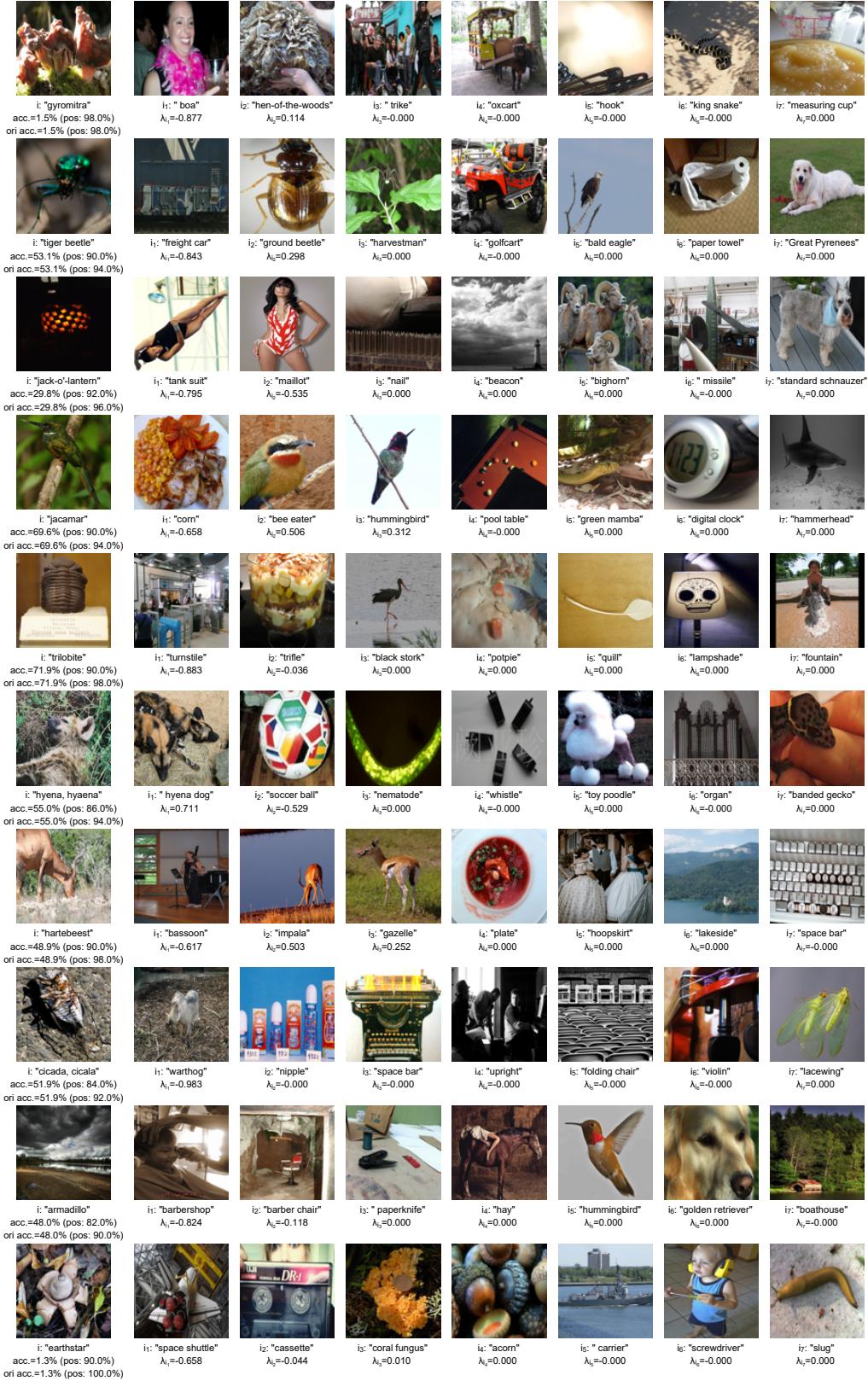


Figure A9: Results from GluMixer-24, where ‘acc.’ and ‘ori acc.’ denote the classification accuracies on the ImageNet validation set, while ‘pos: xx%’ is the accuracy on positive samples only.

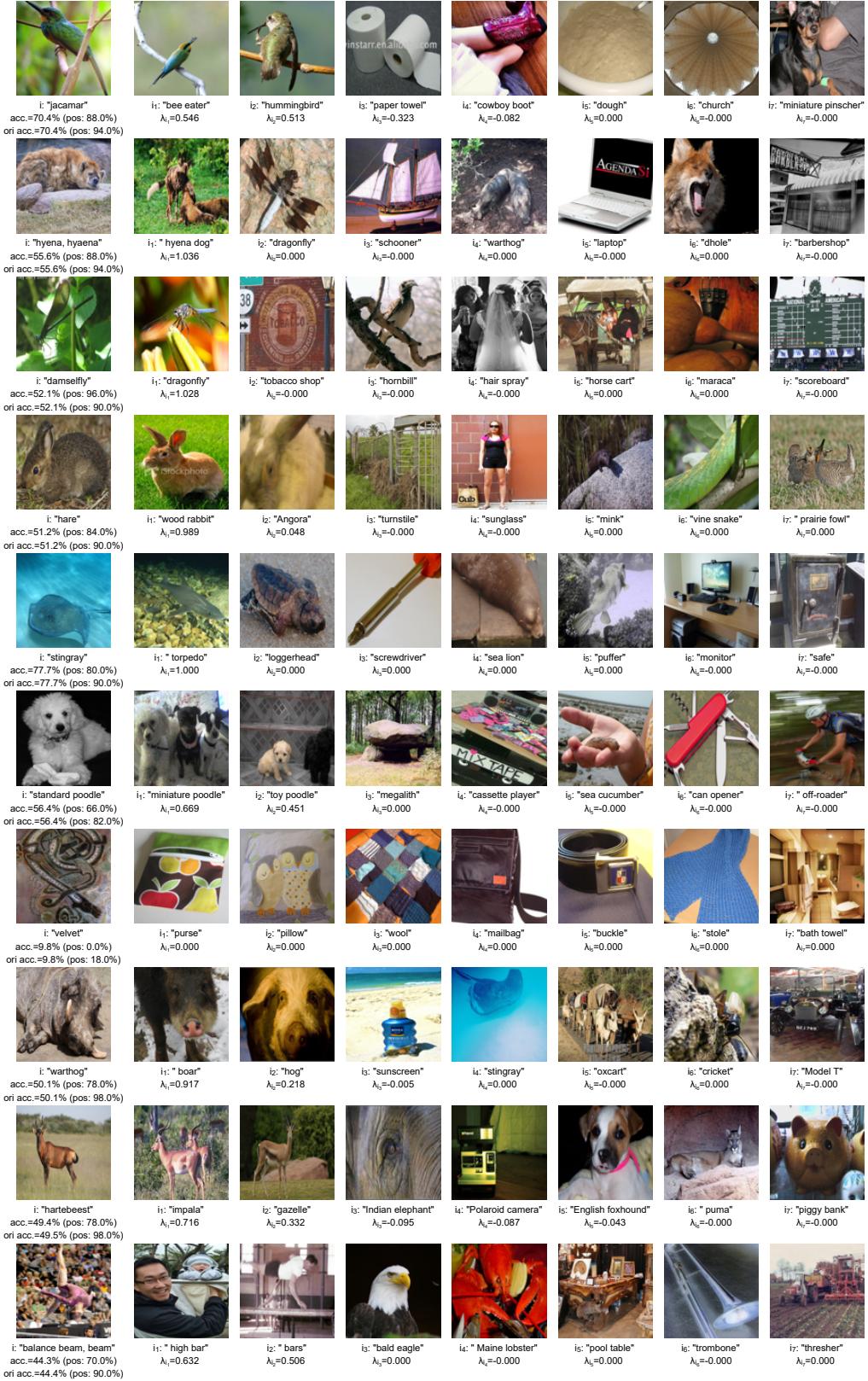


Figure A10: Results from ResMLP-S24, where ‘acc.’ and ‘ori acc.’ denote the classification accuracies on the ImageNet validation set, while ‘pos: xx%’ is the accuracy on positive samples only.

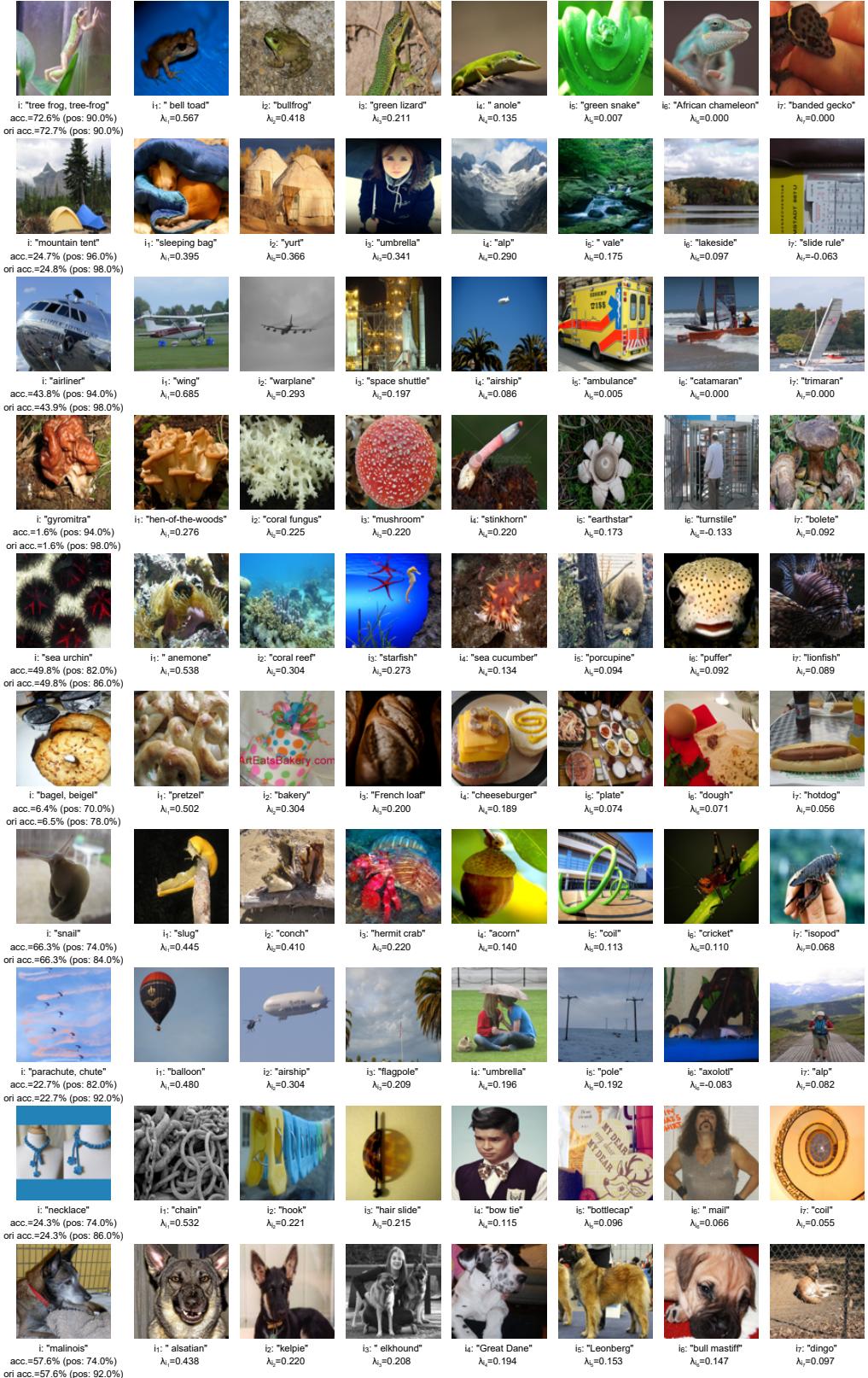


Figure A11: Results from ViT-T, where ‘acc.’ and ‘ori acc.’ denote the classification accuracies on the ImageNet validation set, while ‘pos: xx%’ is the accuracy on positive samples only.

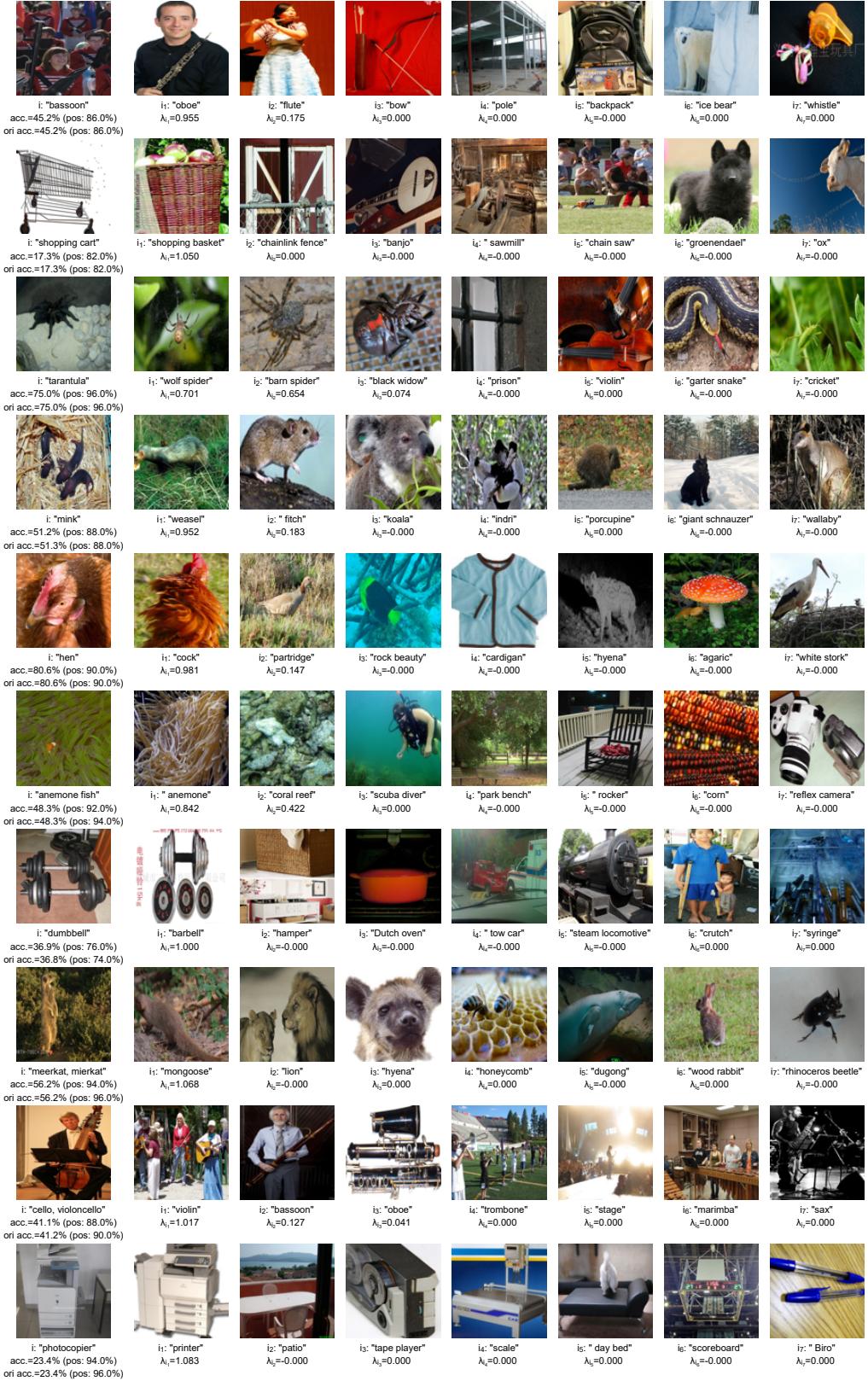


Figure A12: Results from Swin-T, where ‘acc.’ and ‘ori acc.’ denote the classification accuracies on the ImageNet validation set, while ‘pos: xx%’ is the accuracy on positive samples only.

---

**Algorithm A1** Pseudocode of Partial Rank of the Jacobian.

---

```
# image: input images
# model: network
# row_idx, col_idx, patch_size: select a patch of the image to calculate Jacobian matrix

from functools import partial
import torch.nn.functional as functional

def Jacobian_rank(image, model):
    # select a patch of the image to calculate Jacobian matrix
    assert image.size(2) == image.size(3)
    image_size = image.size(2)
    image = image[:, :, row_idx:row_idx + patch_size, col_idx:col_idx + patch_size]
    zero_pad = partial(functional.pad, pad=[(image_size - patch_size) // 2 for _ in range(4)], value=0.)

    # calculate the jacobian matrix
    jacobian_matrix = jacobian(partial(model.net.forward, preprocess=zero_pad), image)

    # adopt trick to predict the singular values
    jacob = jacob.view(-1, image.size(1) * patch_size * patch_size)
    jacob = matmul(jacob.T, jacob)

    # calculate the partial rank of Jacobian matrix
    return matrix_rank(jacob, symmetric=True)
```

---

matmul: matrix multiplication; jacobian: calculate the jacobian matrix; matrix\_rank: calculate the numerical rank of matrix.

---

---

**Algorithm A2** Pseudocode of Perturbed PCA Dimension of Feature Spaces.

---

```
# image: input images
# model: network
# mag_perturb: magnitude of perturbations
# n_perturb: number of perturbations

def Perturbed_Dimension(image, model, mag_perturb=1e-3, n_perturb=5000):
    # extract features with random perturbations
    features = []
    for _ in range(n_perturb):
        # sample random perturbation from Gaussian distribution
        perturb = randn_like(image) * mag_perturb
        # extract feature
        feature = model(image + perturb)
        features.append(feature)
    features = concatenation(features, dim=0)

    # calculate the covariance matrix
    x = input_mean(input, dim=0)
    x = x.view(x.size(0), -1)
    cov_matrix = matmul(x.T, x) # covariance matrix

    # calculate the perturbed PCA dimensions
    return matrix_rank(cov_matrix, symmetric=True)
```

---

matmul: matrix multiplication; randn\_like: sample a random tensor from a Gaussian distribution; matrix\_rank: calculate the numerical rank of matrix.

---

---

**Algorithm A3** Pseudocode of the Classification Dimension of the Final Feature Manifold.

---

```
# image: input images
# model: network
# target: ground-truth labels
# acc_ratio: threshold for measuring intrinsic dimensions of final features

def PCA(X, n_components):
    n = X.shape[0]
    X_mean = mean(X, dim=0, keepdim=True)
    X = X - X_mean
    covariance_matrix = 1 / n * matmul(X.T, X)
    eigenvalues, eigenvectors = evd(covariance_matrix, eigenvectors=True)
    eigenvalues = norm(eigenvalues, dim=1) # modulus of complex numbers
    idx = argsort(-eigenvalues)
    eigenvectors = eigenvectors[:, idx]
    eigenvectors = eigenvectors[:, :n_components]
    return eigenvectors

def Feature_projection(X, V):
    X_proj = zeros_like(X)
    for component_idx in range(V.size(1)):
        eig_vec = V[:, component_idx].unsqueeze(-1)
        eig_vec_norm = eig_vec / norm(eig_vec, p=2, keepdim=True)
        w_proj = matmul(X, eig_vec_norm)
        X_proj_i = w_proj * eig_vec_norm.T
        X_proj += X_proj_i
    return X_proj

def Intrinsic_dimension(image, model, target, acc_ratio=0.95):
    # pre-extract features and calculate original classification accuracy
    feats = model(image) # [n_samples * n_channels]
    acc_ori = calc_acc(feats, target)

    for n_component in range(1, feats.size(1)):
        # compute the eigenvalues and eigenvectors of a real square matrix
        components = PCA(feats, n_component) # [n_channels * n_component]

        # reconstruct features with principal components
        feats_rec = Feature_projection(feats, components)

        # calculate classification accuracy
        acc = calc_acc(feats_rec, target)

        # return classification dimension
        if acc >= acc_ratio * acc_ori:
            return n_component
```

---

matmul: matrix multiplication; evd: eigen value decomposition; calc\_acc: calculating classification accuracy.

---

**Algorithm A4** Pseudocode of Independence Deficit.

---

```
# image: input images
# model: network
# target2index: dictionary mapping from category index to sample indices
# lr: learning rate for Lasso optimization
# n_iteration: number of iterations for Lasso optimization
# w_reg: weight of the L1 regularization term

def Feature_split(feats, class_i, target2index):
    sample_indices = target2index[class_i]
    start_idx, end_idx = sample_indices[0], sample_indices[-1]
    feats_i = feats[start_idx:end_idx+1, :]
    feats_i_n = concatenation((feats_i[:, :class_i], feats_i[:, class_i+1:]), dim=1)
    feats_i_p = feats_i[:, class_i:class_i+1]
    return feats_i_n, feats_i_p

def Independence_deficit(image, model, target2index, lr=1e-5, n_iteration=5000, w_reg=20.0):
    # pre-extract logits
    logits = model(image) # [n_samples * n_classes]

    # Lasso optimization
    for class_i in range(logits.size(1)):
        # split features by category index
        feats_n, feats_p = Feature_split(logits, class_i, target2index)

        # initialize the linear coefficients of category i
        param = Parameter(zeros(feats_n.size(1), 1))

        # start training
        for _ in range(n_iteration):
            loss = mse(matmul(feat_n, param), feat_p) + w_reg * l1_norm(param)
            loss.backward()
            param -= lr * param.grad

        # save trained coefficients
        save(param)
```

---

matmul: matrix multiplication; mse: mean squared error; l1\_norm: sum of the magnitudes of the vectors in a space.