# DECOUPLING DYNAMICAL RICHNESS FROM REPRESENTATION LEARNING: TOWARDS PRACTICAL MEASUREMENT

# **Anonymous authors**

Paper under double-blind review

#### **ABSTRACT**

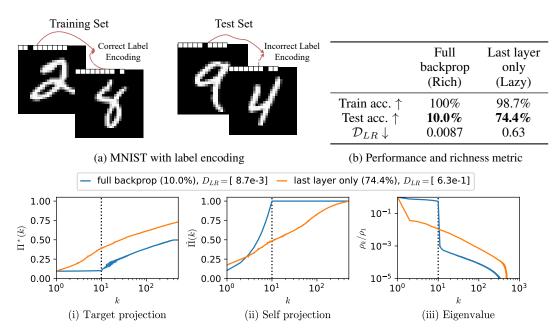
Dynamic feature transformation (the rich regime) does not always align with predictive performance (better representation), yet accuracy is often used as a proxy for richness, limiting analysis of their relationship. We propose a computationally efficient, performance-independent metric of richness grounded in the low-rank bias of rich dynamics, which recovers neural collapse as a special case. The metric is empirically more stable than existing alternatives and captures known lazy-to-rich transitions (e.g., grokking) without relying on accuracy. We further use it to examine how training factors (e.g., learning rate) relate to richness, confirming recognized assumptions and highlighting new observations (e.g., batch normalization promote rich dynamics). An eigendecomposition-based visualization is also introduced to support interpretability, together providing a diagnostic tool for studying the relationship between training factors, dynamics, and representations.

#### 1 Introduction

In machine learning, feature learning is often viewed through two complementary perspectives: improvement of representations and non-linear training dynamics. The *representation perspective* emphasizes feature quality — how well it supports downstream tasks like classification and promotes generalization (Bengio et al., 2013). The *dynamics perspective* — often called rich regime in *rich versus lazy training* (Chizat et al., 2019) — concerns the dynamic transformation of features beyond linear models. While dynamical richness frequently correlates with representation usefulness, rich dynamics reflect a preference (inductive bias) toward certain solutions, without necessarily benefiting all tasks (Geiger et al., 2020; Göring et al., 2025). For example, dynamic feature learning can impair the performance on an image classification task (Fig. 1).

The representation perspective of feature learning — central to deep learning's success — remains poorly understood. However, its dynamic aspects, often shaped in practice by initialization and optimization, are well understood in the context of layerwise linear models (e.g., linear networks). Greedy dynamics with low-rank biases (Saxe et al., 2014; Mixon et al., 2020) explain neural collapse (Papyan et al., 2020) in vision tasks; delayed saturation dynamics (Nam et al., 2024) have been linked to emergence (Brown et al., 2020) and scaling laws (Kaplan et al., 2020) in large language models; and initialization-dependent dynamics (Kumar et al., 2024; Kunin et al., 2024; Lyu et al., 2024) offer insights into grokking (Power et al., 2022) in arithmetic tasks. See Nam et al. (2025) for a comprehensive overview.

To better understand the link between rich dynamics and representation improvement, we need independent metrics for dynamics and representations. However, prior dynamics metrics, while insightful, are not optimized for practical richness measurement (Section 4.1). In this paper, we propose a dynamical richness that compares the activations before and after the last layer. The metric (1) generalizes neural collapse, (2) is computationally cheap, (3) performance-independent, and (4) operates in function space. We demonstrate that our metric empirically tracks dynamical richness without referencing performance, and assess its robustness against existing metrics. Finally, we experiment our metric, with complementary visualization, across various setups, highlighting its potential as a tool for uncovering new empirical insights.



(c) Visualization of representation quality (i) and our metric (ii,iii) using the eigendecomposition

Figure 1: **Rich dynamics**  $\neq$  **better generalization.** We trained a 4-layer MLP on label-encoded MNIST. (a): The first 10 pixels are encoded with true labels in training and random labels in testing; both the encoding and the image serve as valid features for training. (b): A full backpropagation model (*rich*) biases toward the encodings and generalizes poorly, while a last-layer-only-trained model (*lazy*) relies on the full image and generalizes better. Our low-rank-based metric  $\mathcal{D}_{LR}$  (Eq. (4)) quantifies the dynamical richness ( $\mathcal{D}_{LR} \in [0,1]$  where 0 is richest) independent of the performance. (c): Complementary visualization method (Eq. (5)): (i) Cumulative contribution of last-layer features in expressing the target function — top 10 features are irrelevant in the rich model. (ii) Contribution to the learned function — the rich model uses only the top 10 features, while the lazy model uses all. (iii) Relative feature norms — rich model concentrates on the top 10; lazy model decays more gradually. Test accuracies and  $\mathcal{D}_{LR}$  values are shown in parentheses and square brackets, respectively. See Section 5.1 for details, and Appendix C for background and motivation.

# **Contributions:**

- We introduce the Minimum projection operator  $\mathcal{T}_{MP}$  (Definition 1) and define the dynamical low-rank measure  $\mathcal{D}_{LR}$  (Eq. (4)), a light weight and performance independent metric.
- We show that  $\mathcal{D}_{LR}$  reduces to neural collapse as a special case, connecting our formulation to a well-studied phenomenon while extending it to settings without labels.
- We empirically confirm that  $\mathcal{D}_{LR}$  captures known lazy-to-rich transitions such as grokking and target downscaling (Table 2), aligning with theoretical expectations while giving a direct evaluation without referencing performance. We also empirically assess its robustness against alternatives.
- We demonstrate the utility of our metric by showing how various training factors (e.g., architecture and learning rate) relate to richness and performance (Table 3). Our metric explicitly confirms prior assumptions and uncovers new observations, such as batch normalization shifting VGG-16 on CIFAR-100 from lazy to rich dynamics. We further adapt eigendecomposition to visualize the metric for improved interpretability.

## 2 SET UP AND BACKGROUND

For analytical tractability, we use MSE loss and supervised settings where target function entries are orthogonal and isotropic (e.g., balanced C-way classification or regression with scalar output).

**Notations.** For the input space  $\mathcal{X}$ , we train on n samples from q, the underlying true distribution on the input space  $\mathcal{X}$ . The width of the last layer (or the linear model) is p. The dimension of the output function (and the output of the neural network) is C. Denote  $\hat{f}$  and  $f^*$  as the learned function (at any point in training) and the target function, respectively. We use bra-ket notation  $\langle f|g\rangle:=\mathbf{E}_{x\sim q}[f(x)g(x)]$  to mean the inner product over the input distribution (Appendix B). See Appendix A for a **glossary**.

**Empirical details.** Full experimental details, link to the source code, and statistical significance for the tables are provided in Appendix H. All error bars represent one standard deviation.

**Feature kernel operator.** We define the *feature map*  $\Phi: \mathcal{X} \to \mathbb{R}^p$  as the map from the input to the post-activations of the penultimate layer. For  $1 \le k \le p$ , denote the  $k^{\mathrm{th}}$  feature as the  $k^{\mathrm{th}}$  entry  $\Phi_k(x)$  of  $\Phi(x) = [\Phi_1(x), \dots, \Phi_p(x)]$ . We define the *feature kernel operator*  $\mathcal{T}: L^2(\mathcal{X}) \to L^2(\mathcal{X})$ , generalization of self-correlation matrix in function space, as

$$\mathcal{T}[f](x') := \mathbf{E}_{x \sim q} \left[ \sum_{k=1}^{p} \Phi_k(x) \Phi_k(x') f(x) \right], \qquad \mathcal{T} = \sum_{k=1}^{p} |\Phi_k\rangle \langle \Phi_k|, \tag{1}$$

where  $L^2(\mathcal{X})$  is the Hilbert space of square-integrable functions on  $\mathcal{X}$  (with distribution q). The second expression is  $\mathcal{T}$  in bra-ket notation. Using Mercer's theorem (Mercer, 1909), we can decompose the operator into eigenfunctions and eigenvalues

$$\mathcal{T} := \sum_{k=1}^{p} \rho_k |e_k\rangle \langle e_k|, \qquad \mathcal{T}[e_k] = \rho_k e_k, \qquad \langle e_k | e_l\rangle = \delta_{kl}, \tag{2}$$

where  $\rho_k \in \mathbb{R}$  is a non-negative eigenvalue and  $e_k : \mathcal{X} \to \mathbb{R}$  is the orthonormal eigenfunction. The operator  $\mathcal{T}$  and the eigenfunctions  $e_k$  play a key role in kernel regression with feature map  $\Phi$  (i.e.  $K(x,x') = \Phi(x)^T \Phi(x')$ ). The linear model (kernel) has inductive bias toward larger (eigenvalue) eigenfunctions (Bordelon et al., 2020; Jacot et al., 2020; Spigler et al., 2020; Canatar et al., 2021; Cui et al., 2021; El Harzli et al., 2024) and generalizes better when larger eigenfunctions (larger  $\rho_k$ ) also better describe the target function (larger  $\langle e_k | f^* \rangle$ ) — also known as the task-model alignment (Bordelon et al., 2020). See Appendix C for a gentle introduction.

**Rich dynamics.** Chizat et al. (2019) introduced rich dynamics as deviations from the exponential saturation observed in linear models. Subsequent works (Kunin et al., 2024; Dominé et al., 2025; Nam et al., 2025) identified these dynamics as sigmoidal saturation or amplifying dynamics — arising from gradient descent in layerwise models — and connected them to phenomena like neural collapse (Mixon et al., 2020), feature emergence (Nam et al., 2024), and grokking (Kunin et al., 2024). See Nam et al. (2025) for a comprehensive overview.

**Low-rank bias and neural collapse.** Low-rank hidden representations naturally emerge in the rich regime, as shown by studies across linear networks (Saxe et al., 2014; Ji & Telgarsky, 2019; Arora et al., 2019a; Lampinen & Ganguli, 2019; Gidel et al., 2019; Tarmoun et al., 2021), unconstrained feature models (Mixon et al., 2020; Fang et al., 2021), and matrix factorization (Arora et al., 2019b; Li et al., 2020). These works find that gradient dynamics decouple into a minimal number of modes, governed by the rank of the input-output correlation matrix (often that of the output C), causing gradients to concentrate on those modes and produce low-rank representations. In practice, the low-rank bias is observed as neural collapse (Papyan et al., 2020), which consists of four conditions NC1-4 stated in Appendix D. Although often associated with improved representations, neural collapse does not consistently imply better generalization (Zhu et al., 2021; Nguyen et al., 2022; Hui et al., 2022; Su et al., 2023), remaining largely a dynamical phenomenon. See Kothapalli et al. (2022) for an overview.

Change in NTK as measure of rich dynamics. Beyond performance, changes in the Neural tangent kernel (NTK, Jacot et al. (2018)) before and after training are among the most commonly used proxies for rich dynamics (Chizat et al., 2019; Atanasov et al., 2021; Kumar et al., 2024). Although theoretically well-founded, NTK-based measures scale with the total number of parameters and are computationally infeasible even for moderate vision models. For this reason, we omit NTK from our analysis and we focus only on the feature kernel of the last layer (Eq. (1)).

# 3 RICHNESS MEASURE $\mathcal{D}_{LR}$

We introduce a dynamical richness measure that exploits the low-rank bias of rich dynamics and show that it reduces to neural collapse as a special case.

#### 3.1 Low rank bias as richness measure

We define the learned function space  $\hat{\mathcal{H}} = span\{\hat{f}_1,\ldots,\hat{f}_C\}$  where  $\hat{f}_k:\mathcal{X}\to\mathbb{R}$  is the  $k^{\mathrm{th}}$  entry of the network's learned function (at any given time). In an ideal rich dynamics scenario, only the minimal number of features are learned throughout the dynamics and are sufficient to express (linearly span) the learned function space. Leveraging this idea, we define the minimal projection operator  $\mathcal{T}_{MP}$ .

**Definition 1** (Minimum Projection (MP) operator). For a neural network with learned function  $\hat{f}$  and features  $\Phi(x)$ , the corresponding T is an MP-operator  $T_{MP}$  if it can be expressed as

$$\mathcal{T}_{MP}[u] = a_1 \langle \mathbf{1} | u \rangle \mathbf{1} + a_2 P_{\hat{\mathcal{H}}}(u) \quad \text{for all } u \in L^2(\mathcal{X}), \tag{3}$$

where  $a_1, a_2 > 0$ , **1** is a constant function, and  $P_{\hat{\mathcal{H}}}$  is the orthogonal projection onto  $\hat{\mathcal{H}}$ .

Ignoring the constant function (setting  $a_1=0$ ), whose discussion is deferred to Appendix D, the  $\mathcal{T}_{MP}$  is (up to a constant scale) a projection operator  $P_{\hat{\mathcal{H}}}$  that removes all components orthogonal to  $\hat{\mathcal{H}}$ . If  $\mathcal{T}$  is  $\mathcal{T}_{MP}$ , the last-layer features span only a C-dimensional space that matches the learned function space, reflecting the low-rank structure characteristic of rich dynamics. We thus define the low rank measure  $\mathcal{D}_{LR}$  as the similarity between the  $\mathcal{T}$ —defined by the current features—and the MP-operator  $\mathcal{T}_{MP}$ —defined by the current learned function—as a metric for richness:

$$\mathcal{D}_{LR} := 1 - CKA(\mathcal{T}, \mathcal{T}_{MP}), \tag{4}$$

where CKA is the centered kernel alignment (Kornblith et al., 2019) with bounded value in [0,1]. Because CKA is normalized and uses centered (zero-mean) alignment, the metric remains consistent for any  $a_1$  and  $a_2$  (Appendix B.3). We subtract the CKA measure from 1 so that lower values indicate richer dynamics, consistent with more widely used metrics of richness. The novelty of our metric  $D_{LR}$  lies in defining the minimum projection operator  $\mathcal{T}_{MP}$  and comparing it to T to **quantify dynamical richness**, not in using well-established CKA. Notably, CKA is primarily used to compare the NTK before and after training (Chizat et al., 2019; Kumar et al., 2024; Baratin et al., 2021).

**Intuition.** Our metric (Eq. (4)) compares the activations **before** ( $\mathcal{T}$ ) and **after** ( $\mathcal{T}_{MP}$ ) the last layer. In rich dynamics, low-rank bias dictates that only the minimal necessary features are learned and used ( $\mathcal{D}_{LR}=0$ ). In this case, the final layer performs no additional processing since earlier layers have already completed the task. Conversely, an excess of features indicates a bottleneck — such as limited expressivity in earlier layers — that prevents the full manifestation of low-rank bias.

**Time complexity.** The algorithm is highly efficient. For n samples (from either train or test), it requires n forward passes to record activations before and after the final layer, producing matrices of size  $n \times p$  and  $n \times C$ , where p is the **last-layer width** and C the number of classes.  $\mathcal{D}_{LR}$  is then computed in  $\mathcal{O}(npC)$ . Since typically  $n \gg p \gg C$ , setting  $n \approx \mathcal{O}(p)$  suffices, reducing complexity to  $\mathcal{O}(p^2C)$ . In standard models with  $p \approx 10^3$ , this is far cheaper than NTK-based methods, which scale quadratically with the **total number of parameters**. See Appendix E for details.

#### 3.2 Connection to Neural Collapse

Suppose that the empirical distribution coincides with the true distribution and the neural network perfectly classifies all labels with one-hot vectors. We show that if  $\mathcal{T}$  is an MP-operator, then the NC1 and NC2 conditions (Appendix D) of neural collapse hold. Let dagger  $\dagger$  denote pseudo inverse,  $A_i$  the set of datapoints with class label i, so the  $i^{\text{th}}$  learned function  $\hat{f}_i$  coincides with the one-hot indicator function for  $A_i$ . Following Papyan et al. (2020), we define the  $i^{\text{th}}$  class mean vector as  $\mu_i := \mathbf{E}_{x \in A_i}[\Phi(x)]$  and the global mean vector as  $\bar{\mu} := C^{-1} \sum_{i=1}^{C} \mu_i$ .

 **Proposition 1.** If  $\mathcal{T}$  is an MP-operator, then the NC1 condition (collapse of within-class variability)  $\Sigma_b^I \Sigma_W = 0$  holds, where inter-class covariance matrix  $\Sigma_b$  and intra-class covariance matrix  $\Sigma_W$  are

$$\Sigma_b = \sum_{j=1}^{C} (\mu_j - \bar{\mu})(\mu_j - \bar{\mu})^T \quad \text{and} \quad \Sigma_W = \sum_{i=1}^{C} \mathbf{E}_{x \in A_i} \left[ (\Phi(x) - \mu_i)(\Phi(x) - \mu_i)^T \right].$$

**Proposition 2.** If T is an MP-operator, then the NC2 condition (convergence of features to a simplex equiangular tight frame) holds:

$$(\mu_i - \bar{\mu})^T (\mu_j - \bar{\mu}) \propto \delta_{ij} - \frac{1}{C}.$$

See Appendix D for proofs and further discussions.

Although the ideal criteria of  $\mathcal{T}$  is  $\mathcal{T}_{MP}$  implies neural collapse criteria as a special case, they differ in general: we measure how well features (random variables or functions) express the learned function, while neural collapse concerns how training feature vectors represent class mean vectors. Based on function space, our measure extends to test data, enables feature quality assessment (Eq. (5)), and is more empirically robust (Tables 1 and 2), extending beyond neural collapse (see Appendix D).

### EXPERIMENTS

Here, we share our empirical results. We first compare our metric with existing metrics of richness, then we confirm that our measure empirically tracks known lazy-to-rich transitions. We further conduct experiments on various training setup to explore the relationship between training factor, rich dynamics, and performance (better representations).

#### 4.1 COMPARISON TO PRIOR MEASURES OF RICH DYNAMICS

Three alternative metrics for richness measures are: (1) similarity to the initial kernel,  $S_{init}$  =  $CKA(K_{\text{init}}, K_{\text{learned}}) \in [0, 1]$  (Yang & Hu, 2021), (2) the parameter norm  $\|\theta\|_F^2$  (Lyu et al., 2024), and (3) class separation from neural collapse,  $NC1 = Tr(\Sigma_h^{\dagger} \Sigma_W)$  (Papyan et al., 2020; Stevens et al., 2002; He & Su, 2023; Xu & Liu, 2023; Súkeník et al., 2024), where  $\Sigma_h$  and  $\Sigma_W$  are the interand intra-class covariances (Proposition 1). Prior metrics depend on the initial kernel, parameter norms, or class labels, which can constrain their use as independent measures of richness.

Table 1: Richness for weight decay

Table 2: Richness metrics for target downscaling

Epoch	0 (init)	200	$\alpha$	$2\cdot 10^{-1}$	$2 \cdot 10^{0}$	$2 \cdot 10^1$
Train Acc.↑	10.0%	10.2%	Train Acc.↑	100 %	100%	100%
Test Acc.↑	9.85%	10.0%	Test Acc.↑	92.7%	92.4%	88.3%
${\cal D}_{LR}\!\!\downarrow$	0.59	1.0	${\cal D}_{LR}\!\!\downarrow$	$4.9\cdot10^{-2}$	$1.1\cdot 10^{-1}$	$5.6\cdot10^{-1}$
$\mathcal{S}_{init} \!\!\downarrow$	1.0	0.20	$\mathcal{S}_{init} \!\!\downarrow$	$6.8 \cdot 10^{-2}$	$4.1 \cdot 10^{-2}$	$5.2 \cdot 10^{-2}$
$\ \theta\ _F^2\downarrow$	$3.1 \cdot 10^{3}$	$2.2 \cdot 10^{-5}$	$\ \theta\ _F^2 \downarrow$	$3.4\cdot 10^3$	$3.2\cdot 10^3$	$3.2 \cdot 10^{3}$
$NC1\downarrow$	$1.2 \cdot 10^{5}$	$7.5 \cdot 10^{-14}$	${N}{C}1\downarrow$	$2.3 \cdot 10^{4}$	$3.2 \cdot 10^{3}$	$8.1 \cdot 10^{2}$

Table 1 shows an extreme case of training an MLP on MNIST with large weight decay and a negligible learning rate. Here, the dynamics are dominated by L2 weight decay, with little meaningful learning. While existing metrics can sometimes misinterpret this as rich behavior (smaller after training), our span-based measure correctly identifies the lack of dynamical richness (bigger after training).

Table 2 presents a more practical setting, where we tune laziness via target downscaling (i.e.,  $y \to y/\alpha$ ). Prior works (Chizat et al., 2019; Geiger et al., 2020) show that scaling targets by a factor  $\alpha$  induces lazier training where larger  $\alpha$  implies greater laziness. A good metric should capture this. Our measure aligns with  $\alpha$ , while all other measures misalign with laziness. By not relying on initial kernel, weight norm, or labels, our metric shows greater robustness in this setup. See Fig. 11 in Appendix F for visualization of the dynamics in Table 2.

**Kernel distance from initialization** ( $S_{init}$ ). Kernel deviation (Yang & Hu, 2021), albeit often the comparison of NTKs, is a common measure of rich dynamics (Chizat et al., 2019). However, it measures the deviation from *initial kernel*, not *how* the kernel changed (e.g. weight decay alone changes the kernel in Table 1). While theoretically appealing for near-lazy models, the metric can be less accurate in deep, rich-training regimes where most practical training occurs.

**Parameter norm** ( $\|\theta\|_F^2$ ). Smaller parameter norms often correlate with rich dynamics in practice, similar to how the authors of Lyu et al. (2024) used it in their study of grokking. However, small weights promote rich dynamics (Kumar et al., 2024; Atanasov et al., 2024; Nam et al., 2025; Saxe et al., 2014; 2019), not the other way around. In fact, rich dynamics can occur for larger initialization (Braun et al., 2022; Dominé et al., 2025) and depends on broader factors such as layer imbalances (Kunin et al., 2024; Dominé et al., 2025; Nam et al., 2025).

**Neural collapse measure** (NC1). Separation-based metrics measure how training samples deviate relative to class boundaries (Papyan et al., 2020; Stevens et al., 2002; He & Su, 2023; Xu & Liu, 2023; Súkeník et al., 2024). While most similar to our metric, they are unbounded, sensitive to output scaling, and can be empirically ill-conditioned. As shown in Table 1, these issues can lead to dramatic value shifts and hinder interpretability.

**Low-rank measure**  $(\mathcal{D}_{LR})$ . Our proposed metric evaluates the alignment between features and the learned function, achieving its optimum when they span the same space isotropically (i.e., a scaled projection operator). Our metric exploits the low-rank bias of rich dynamics through the alignment and is normalized between [0,1]. Crucially, it does so without relying on class labels, accuracy, or the initial kernel, making the measure more appealing as an independent measure of richness.

#### 4.2 Empirical findings regarding training factors

Table 3: Richness measure and performance on various setups

Task	Architecture	Condition	Train acc.	Test acc.	$D_{LR}$	Figure
Mod	2-layer	Step 200 (before grokking)	100%	5.2%	0.51	Fig. 12
97 transformer		Step 3000 (after grokking)	100%	99.8%	0.11	Fig. 12
CIFAR-		learning rate = 0.005	100%	66.3%	0.053	
100	ResNet18	learning rate = $0.05$	100%	78.3%	0.025	Fig. 4
100		learning rate = $0.2$	100%	74.5%	0.039	
CIFAR-	ResNet18	weight decay = 0	100%	93.5%	0.05	Fig. 13
10	RESINCTIO	weight decay = $10^{-4}$	100%	94.1%	0.015	11g. 13
		weight decay = $10^{-3}$	100%	94.8%	0.003	
CIFAR-	ResNet18	This experiment compares	100%	94.8%	0.026	Fig. 14
10	MLP	architectures only.	99.8%	55.4%	0.48	Fig. 14
CIFAR-		no label shuffling	100%	95.0%	0.031	
10	ResNet18	10% label shuffling	100%	66.1%	0.042	Fig. 15
10		full label shuffling	100%	9.5%	0.034	
MNIST	CNN	full backpropagation	100%	99.1%	0.043	Fig. 16
MINIST	CIVIN	last layer only training	99.7%	96.8%	0.51	11g. 10
CIFAR-	VGG-16	without batch nomralization	99.5%	21.7%	0.66	Fig. 3
100	V UU-10	with batch normalization	100%	72.0%	0.073	1 1g. 3
		· · · · · · · · · · · · · · · · · · ·				

Table 3 demonstrates the practical usefulness of our performance-independent metric, summarizing the correlation between training factor, performance, and richness in various setups. The visualization of most experiments are provided in Appendix F.

The first row empirically confirms that our metric is a richness metric, explicitly capturing known lazy-to-rich transition of grokking (Kumar et al., 2024; Kunin et al., 2024; Lyu et al., 2024; Nam et al., 2025).

Rows two to four confirm the assumptions that optimal learning rate ( $2^{nd}$  row), weight decay ( $3^{rd}$  row), and architecture ( $4^{th}$  row) achieves high performance **through rich dynamics** (Ginsburg et al.,

2018; Liu et al., 2022; He et al., 2016). Indeed, our metric makes this link **explicit** by having a smallest value for the optimal setting. Although such relationships are broadly recognized, they are typically only addressed implicitly.

 The fifth and sixth rows, like Fig. 1, demonstrate that rich dynamics do not strictly correlate with performance. ResNet18 on CIFAR-10 retains rich dynamics even when the labels are shuffled, and convolutional neural network (CNN) on MNIST can achieve similar performance through both lazy and rich dynamics.

The last row reveals a **new observation**: VGG-16 on CIFAR-100 is lazy without batch normalization but rich with it, accompanied by a significant performance gap. While the performance effect of batch normalization is empirically established, its underlying role remains debated. Our metric helps clarify this by reframing generalization in terms of the more tractable problem of rich dynamics.

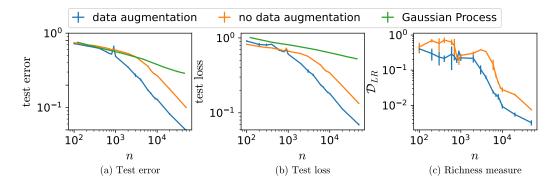


Figure 2: Learning curve and feature learning metric. (a): Learning curves of ResNet18 on CIFAR-10. Both error (a) and loss (b) learning curves show a transition to a faster-decaying power law with additional data near  $n \approx 10^3$ , correlating with the shift in decay of the richness measure  $\mathcal{D}_{LR}$  in (c). This agrees with theoretical study on phase transition (Rubin et al., 2024) that a sufficiently large number of data points is critical for rich dynamics — a promising observation toward better understanding feature learning dynamics. A linear model (Gaussian process) was plotted in (a,b) to highlight the transition into faster-decaying learning curve.

#### 5 VISUALIZATION METHODS

While metrics provide quantitative summaries, component-wise visualization aids interpretation. For instance, how significant is the contrast between VGG-16 trained with versus without batch normalization (Fig. 3)? Or why does a small change in  $\mathcal{D}_{LR}$  correspond to large performance differences under varying learning rates (Fig. 4)? Does the CIFAR-10 MLP begin in a lazy regime, or transition into it during training (Fig. 14 in Appendix F)? To address such questions, we introduce a complementary visualization based on widely used eigendecomposition of kernel (Eq. (1)).

#### 5.1 VISUALIZATION THROUGH DECOMPOSED FEATURES

Bengio et al. (2013) described feature learning as the process of learning better representations for a downstream model, such as the classifier. Building on this idea, we view the earlier layers, represented by the feature map  $\Phi: \mathcal{X} \to \mathbb{R}^p$ , as providing improved features for the final linear layer. Combined with our  $\mathcal{T}$ -dependent richness metric, this motivates a visualization method that extends the linear model analysis of  $\mathcal{T}$ . See Appendix C for related works and a gentle overview.

Our visualizes quantifies three aspects (cumulative quality  $\Pi^*(k)$ , cumulative utilization  $\hat{\Pi}(k)$ , and relative eigenvalue  $\rho_k/\rho_1$ ) of the last-layer features through the eigenfunctions  $e_k$  of  $\mathcal{T}$ :

$$(i): \Pi^{*}(k) = \sum_{j=1}^{k} \frac{\langle e_{k} | P_{\mathcal{H}^{*}}[e_{k}] \rangle^{2}}{\dim(\mathcal{H}^{*})}, \qquad (ii): \hat{\Pi}(k) = \sum_{j=1}^{k} \frac{\langle e_{k} | P_{\hat{\mathcal{H}}}[e_{k}] \rangle^{2}}{\dim(\hat{\mathcal{H}})}, \qquad (iii): \rho_{k}/\rho_{1}, \quad (5)$$

where  $\mathcal{H}^* = span\{f_1^*, \dots, f_C^*\}$  and  $\hat{\mathcal{H}} = span\{\hat{f}_1, \dots, \hat{f}_C\}$  are the target and learned function spaces, and  $P_{\mathcal{H}}: \mathcal{H} \to \mathcal{H}$  is the projection operator onto the space  $\mathcal{H}$ . All cumulative measures lie in [0,1], reflecting how well the top k features span the respective space. Notably,  $\Pi^*(k)$  is the cumulative power in Canatar & Pehlevan (2022), quantifying the contribution of the top features in expressing the target.

These three measures together capture complementary views of feature learning. The cumulative quality  $(\Pi^*(k))$  reflects how well the features align with the task. The utilization  $(\hat{\Pi}(k))$  indicates how many features are used by the final layer, while the relative eigenvalues  $(\rho_k/\rho_1)$  show their relative magnitudes or importance. The latter two  $(\hat{\Pi}$  and  $\rho_k/\rho_1)$  decompose the deviation from  $\mathcal{D}_{LR}=0$  condition by visualizing how many features are used and how many are significant (non-negligible).

For computation, we approximate eigenfunctions with the Nyström method (Baker & Taylor, 1979; Williams & Rasmussen, 2006) using a sufficiently large sample size (Appendix E). The function approximation enables test-set evaluation of the quality measure in (i), analogous to using test loss as a proxy for generalization, ultimately visualizing both performance (i) and the metric (ii, iii) in a common basis. As eigendecomposition arises naturally in study of kernels, similar visualizations have appeared in prior works (e.g., (Canatar et al., 2021)); see Appendix G for details.

Interpreting visualization In perfectly rich dynamics (e.g., blue line in Fig. 3),  $\hat{\Pi}(C) = 1$  in (ii), showing that only the first C features are used by the last layer. Additionally,  $\rho_2 \approx \rho_C$  and  $\rho_C \gg \rho_{C+1}$  in (iii), indicating that the features have only C significant dimensions. The first eigenvalue, corresponding to a constant function, is ignored because of the centering in CKA (see Appendix D). If the performance is high,  $\hat{\Pi}(C) \approx 1$  in (i), showing that first C features well express the target. For the lazy regime, more eigenfunction will be used (slower capping in (ii)) and the feature occupy higher dimensions (slower decay in (iii)).

### 5.2 VISUALIZATION RESULTS

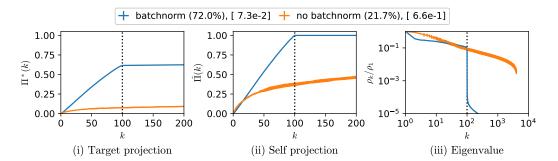


Figure 3: **Visualization of VGG16 on CIFAR-100 with and without batch normalization.** We visualize the last row of Table 3, where batch normalization shifts the model from the lazy to the rich regime. The eigenvalue distribution (iii) highlights this difference: with batch normalization, only 100 features are significant, whereas without it the eigenvalues decay slowly.

Figs. 3 and 4 illustrate how visualization complements the metric. In Fig. 3, the batch-norm included model uses only 100 features (eigenfunctions) in (ii) and shows clear low-rank structure in (iii). The batch normalization lacking model uses significantly more feature (4096 in total) in (ii) and show power-law distribution in (iii): indicating significantly lazy dynamics. In Fig. 4(ii), the model with the smallest learning rate uses over 100 features to express the learned function, unlike models with larger learning rates, suggesting lazier dynamics and poorer performance.

Fig. 5 reveals another novel pattern in feature learning dynamics: feature quality correlates with feature intensity during training, with larger features improving faster. While it is expected that a generalizing model in the rich regime obtains a few high-quality features after training, the correlation between quality and intensity *during* training has not been previously observed or studied. Interestingly, it appears to contrast with the silent alignment (Atanasov et al., 2021), which suggests quality (dynamically) precedes intensity — highlighting a new direction for future research.

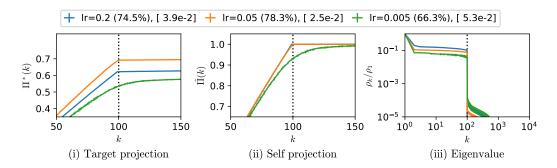


Figure 4: **Visualization on the role of learning rate.** We visualize the  $2^{nd}$  row of Table 3 where the learning rates are varied (up to training instability) for ResNet18 on CIFAR-100. The second column (ii) shows that smallest learning rate uses significantly more eigenfunctions (features), while other models uses minimal 100 eigenfunctions, indicating a lazier dynamics.

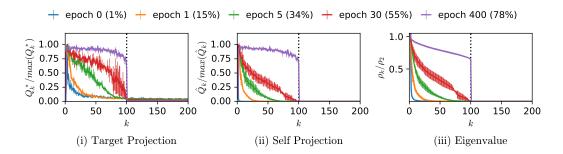


Figure 5: Correlation among dynamics of feature quality, utilization, and intensity. We show individual metrics (e.g.,  $Q^*(k) := \Pi^*(k) - \Pi^*(k-1)$ ) instead of cumulative metrics ( $\Pi^*(k)$  and  $\hat{\Pi}(k)$ ) at different epochs for ResNet18 on CIFAR-100, normalized for better presentation. Larger intensity features exhibit higher quality and utilization during training.

#### 6 DISCUSSION AND CONCLUSION

**Discussion.** We introduced a performance-independent metric for dynamical richness,  $\mathcal{D}_{LR}$ , which reduces to neural collapse as a special case (Propositions 1 and 2) and empirically agrees with prior studies on lazy-to-rich transitions (Tables 2 and 3).

Table 3, Fig. 2 are examples of how the metric can be used for the analysis. We confirm established assumptions explicitly with an independent metric, while potentially challenging a few (Fig. 5). We also reveal new findings (Fig. 3), offering a direction for theoretical research.

Moreover, we provide a complementary visualization to extend these results for additional information, also uncovering novel insights such as the alignment between representation quality and feature intensity (Fig. 5).

**Limitations.** The current form of  $\mathcal{T}_{MP}$  is limited to orthogonal and isotropic target functions. While this covers most classification tasks, a more general setup would be preferable. Additionally, our metric is light-weight, but as a trade-off, focuses only on last-layer features, leaving the dynamics of the entire network dynamics unexplored.

**Conclusion.** By offering a lightweight, robust, and performance-independent metric on rich dynamics, we aim to lay the groundwork for future theoretical studies on their connection to representation learning. More broadly, we see this work as a diagnostic tool toward bridging empirical observations on representation learning (e.g., rule of thumb) and theoretical understanding of rich dynamics. In future work, we plan to extend beyond the current focus on balanced tasks.

# REPRODUCIBILITY STATEMENT

As mentioned in the main text, full experimental details, link to the source code, and statistical significance for the tables are provided in Appendix H. All error bars represent one standard deviation. The algorithm for approximating the metric and visualization are provided in Appendix E.

#### REFERENCES

- Hervé Abdi and Lynne J Williams. Principal component analysis. *Wiley interdisciplinary reviews:* computational statistics, 2(4):433–459, 2010.
- Andrea Amadei, Antonius BM Linssen, and Herman JC Berendsen. Essential dynamics of proteins. *Proteins: Structure, Function, and Bioinformatics*, 17(4):412–425, 1993.
- Sanjeev Arora, Nadav Cohen, Noah Golowich, and Wei Hu. A convergence analysis of gradient descent for deep linear neural networks. In 7th International Conference on Learning Representations, 2019a. URL https://openreview.net/forum?id=SkMQg3C5K7.
- Sanjeev Arora, Nadav Cohen, Wei Hu, and Yuping Luo. Implicit regularization in deep matrix factorization. *Advances in Neural Information Processing Systems*, 32, 2019b.
- Alexander Atanasov, Blake Bordelon, and Cengiz Pehlevan. Neural networks as kernel learners: The silent alignment effect. *arXiv preprint arXiv:2111.00034*, 2021.
- Alexander Atanasov, Alexandru Meterez, James B Simon, and Cengiz Pehlevan. The optimization landscape of sgd across the feature learning strength. *arXiv preprint arXiv:2410.04642*, 2024.
- Christopher TH Baker and RL Taylor. The numerical treatment of integral equations. *Journal of Applied Mechanics*, 46(4):969, 1979.
- Aristide Baratin, Thomas George, César Laurent, R Devon Hjelm, Guillaume Lajoie, Pascal Vincent, and Simon Lacoste-Julien. Implicit regularization via neural feature alignment. In *International Conference on Artificial Intelligence and Statistics*, pp. 2269–2277. PMLR, 2021.
- Yoshua Bengio, Aaron Courville, and Pascal Vincent. Representation learning: A review and new perspectives. *IEEE transactions on pattern analysis and machine intelligence*, 35(8):1798–1828, 2013.
- Blake Bordelon, Abdulkadir Canatar, and Cengiz Pehlevan. Spectrum dependent learning curves in kernel regression and wide neural networks. In *International Conference on Machine Learning*, pp. 1024–1034. PMLR, 2020.
- Lukas Braun, Clémentine Dominé, James Fitzgerald, and Andrew Saxe. Exact learning dynamics of deep linear networks with prior knowledge. *Advances in Neural Information Processing Systems*, 35:6615–6629, 2022.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- Kevin M Buckley and Xiao-Liang Xu. Spatial-spectrum estimation in a location sector. *IEEE transactions on acoustics, speech, and signal processing*, 38(11):1842–1852, 1990.
- Abdulkadir Canatar and Cengiz Pehlevan. A kernel analysis of feature learning in deep neural networks. In 2022 58th Annual Allerton Conference on Communication, Control, and Computing (Allerton), pp. 1–8. IEEE, 2022.
- Abdulkadir Canatar, Blake Bordelon, and Cengiz Pehlevan. Spectral bias and task-model alignment explain generalization in kernel regression and infinitely wide neural networks. *Nature communications*, 12(1):2914, 2021.
  - Lenaic Chizat, Edouard Oyallon, and Francis Bach. On lazy training in differentiable programming. *Advances in neural information processing systems*, 32, 2019.

- Hugo Cui, Bruno Loureiro, Florent Krzakala, and Lenka Zdeborová. Generalization error rates
   in kernel regression: The crossover from the noiseless to noisy regime. Advances in Neural
   Information Processing Systems, 34:10131–10143, 2021.
  - Li Deng. The mnist database of handwritten digit images for machine learning research. *IEEE Signal Processing Magazine*, 29(6):141–142, 2012.
  - Clémentine Carla Juliette Dominé, Nicolas Anguita, Alexandra Maria Proca, Lukas Braun, Daniel Kunin, Pedro A. M. Mediano, and Andrew M Saxe. From lazy to rich: Exact learning dynamics in deep linear networks. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=ZXaocmXc6d.
  - Ouns El Harzli, Bernardo Cuenca Grau, Guillermo Valle-Pérez, and Ard A Louis. Double-descent curves in neural networks: a new perspective using gaussian processes. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pp. 11856–11864, 2024.
  - Cong Fang, Hangfeng He, Qi Long, and Weijie J Su. Layer-peeled model: Toward understanding well-trained deep neural networks. *arXiv* preprint arXiv:2101.12699, 4, 2021.
  - Mario Geiger, Stefano Spigler, Arthur Jacot, and Matthieu Wyart. Disentangling feature and lazy training in deep neural networks. *Journal of Statistical Mechanics: Theory and Experiment*, 2020 (11):113301, 2020.
  - Gauthier Gidel, Francis Bach, and Simon Lacoste-Julien. Implicit regularization of discrete gradient dynamics in linear neural networks. *Advances in Neural Information Processing Systems*, 32, 2019.
  - Boris Ginsburg, Igor Gitman, and Yang You. Large batch training of convolutional networks with layer-wise adaptive rate scaling, 2018. URL https://openreview.net/forum?id=rJ4uaX2aW.
  - Gene H Golub et al. Cf vanloan, matrix computations. *The Johns Hopkins*, 113(10):23–36, 1996.
  - Niclas Alexander Göring, Charles London, Abdurrahman Hadi Erturk, Chris Mingard, Yoonsoo Nam, and Ard A. Louis. Feature learning is decoupled from generalization in high capacity neural networks. In *High-dimensional Learning Dynamics* 2025, 2025. URL https://openreview.net/forum?id=pVJpk4tuVu.
  - Hangfeng He and Weijie J Su. A law of data separation in deep learning. *Proceedings of the National Academy of Sciences*, 120(36):e2221704120, 2023.
    - Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
  - Mark O Hill. Diversity and evenness: a unifying notation and its consequences. *Ecology*, 54(2): 427–432, 1973.
  - Like Hui, Mikhail Belkin, and Preetum Nakkiran. Limitations of neural collapse for understanding generalization in deep learning. *arXiv preprint arXiv:2202.08384*, 2022.
  - Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and generalization in neural networks. In *Advances in neural information processing systems*, pp. 8571–8580, 2018.
  - Arthur Jacot, Berfin Simsek, Francesco Spadaro, Clément Hongler, and Franck Gabriel. Kernel alignment risk estimator: Risk prediction from training data. *Advances in Neural Information Processing Systems*, 33:15568–15578, 2020.
  - Ziwei Ji and Matus Telgarsky. Gradient descent aligns the layers of deep linear networks. In *International Conference on Learning Representations*, 2019. URL https://openreview.net/forum?id=HJflg30qKX.

- Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott
   Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models.
   arXiv preprint:2001.08361, 2020.
  - Simon Kornblith, Mohammad Norouzi, Honglak Lee, and Geoffrey Hinton. Similarity of neural network representations revisited. In *International conference on machine learning*, pp. 3519–3529. PMLR, 2019.
  - Vignesh Kothapalli, Ebrahim Rasromani, and Vasudev Awatramani. Neural collapse: A review on modelling principles and generalization. *arXiv preprint arXiv:2206.04041*, 2022.
  - Alex Krizhevsky, Vinod Nair, and Geoffrey Hinton. Cifar-10 (canadian institute for advanced research). 2009. URL http://www.cs.toronto.edu/~kriz/cifar.html.
  - Tanishq Kumar, Blake Bordelon, Samuel J. Gershman, and Cengiz Pehlevan. Grokking as the transition from lazy to rich training dynamics. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=vt5mnLVIVo.
  - Daniel Kunin, Allan Raventos, Clémentine Carla Juliette Dominé, Feng Chen, David Klindt, Andrew M Saxe, and Surya Ganguli. Get rich quick: exact solutions reveal how unbalanced initializations promote rapid feature learning. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL https://openreview.net/forum?id=eNM94i7R3A.
  - Andrew K. Lampinen and Surya Ganguli. An analytic theory of generalization dynamics and transfer learning in deep linear networks. In 7th International Conference on Learning Representations, 2019. URL https://openreview.net/forum?id=ryfMLoCqtQ.
  - Zhiyuan Li, Yuping Luo, and Kaifeng Lyu. Towards resolving the implicit bias of gradient descent for matrix factorization: Greedy low-rank learning. *arXiv* preprint arXiv:2012.09839, 2020.
  - Ziming Liu, Ouail Kitouni, Niklas S Nolte, Eric Michaud, Max Tegmark, and Mike Williams. Towards understanding grokking: An effective theory of representation learning. *Advances in Neural Information Processing Systems*, 35:34651–34663, 2022.
  - Yizhang Lou, Chris E Mingard, and Soufiane Hayou. Feature learning and signal propagation in deep neural networks. In *International Conference on Machine Learning*, pp. 14248–14282. PMLR, 2022.
  - Kaifeng Lyu, Jikai Jin, Zhiyuan Li, Simon Shaolei Du, Jason D. Lee, and Wei Hu. Dichotomy of early and late phase implicit biases can provably induce grokking. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=XsHqr9dEGH.
  - James Mercer. Xvi. functions of positive and negative type, and their connection the theory of integral equations. *Philosophical transactions of the royal society of London. Series A, containing papers of a mathematical or physical character*, 209(441-458):415–446, 1909.
  - Dustin G Mixon, Hans Parshall, and Jianzong Pi. Neural collapse with unconstrained features. *arXiv* preprint arXiv:2011.11619, 2020.
  - Ari Morcos, Maithra Raghu, and Samy Bengio. Insights on representational similarity in neural networks with canonical correlation. *Advances in neural information processing systems*, 31, 2018.
  - Preetum Nakkiran, Gal Kaplun, Yamini Bansal, Tristan Yang, Boaz Barak, and Ilya Sutskever. Deep double descent: Where bigger models and more data hurt. *Journal of Statistical Mechanics: Theory and Experiment*, 2021(12):124003, 2021.
  - Yoonsoo Nam, Nayara Fonseca, Seok Hyeong Lee, Chris Mingard, and Ard A. Louis. An exactly solvable model for emergence and scaling laws in the multitask sparse parity problem. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL https://openreview.net/forum?id=cuWsR25bbI.

Yoonsoo Nam, Seok Hyeong Lee, Clémentine Carla Juliette Dominé, Yeachan Park, Charles London, Wonyl Choi, Niclas Alexander Göring, and Seungjai Lee. Position: Solve layerwise linear models first to understand neural dynamical phenomena (neural collapse, emergence, lazy/rich regime, and grokking). In *Forty-second International Conference on Machine Learning Position Paper Track*, 2025. URL https://openreview.net/forum?id=nrlGUdlo16.

- Duc Anh Nguyen, Ron Levie, Julian Lienen, Gitta Kutyniok, and Eyke Hüllermeier. Memorization-dilation: Modeling neural collapse under label noise. *arXiv preprint arXiv:2206.05530*, 2022.
- Vardan Papyan. Traces of class/cross-class structure pervade deep learning spectra. *Journal of Machine Learning Research*, 21(252):1–64, 2020. URL http://jmlr.org/papers/v21/20-933.html.
- Vardan Papyan, XY Han, and David L Donoho. Prevalence of neural collapse during the terminal phase of deep learning training. *Proceedings of the National Academy of Sciences*, 117(40): 24652–24663, 2020.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems 32*, pp. 8024–8035. Curran Associates, Inc., 2019. URL http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf.
- Alethea Power, Yuri Burda, Harri Edwards, Igor Babuschkin, and Vedant Misra. Grokking: Generalization beyond overfitting on small algorithmic datasets. *arXiv:2201.02177*, 2022.
- Maithra Raghu, Justin Gilmer, Jason Yosinski, and Jascha Sohl-Dickstein. Svcca: Singular vector canonical correlation analysis for deep learning dynamics and interpretability. *Advances in neural information processing systems*, 30, 2017.
- Noa Rubin, Inbar Seroussi, and Zohar Ringel. Grokking as a first order phase transition in two layer networks. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=3ROGsTX3IR.
- Andrew Saxe, Shagun Sodhani, and Sam Jay Lewallen. The neural race reduction: Dynamics of abstraction in gated networks. In *International Conference on Machine Learning*, pp. 19287–19309. PMLR, 2022.
- Andrew M. Saxe, James L. McClelland, and Surya Ganguli. Exact solutions to the nonlinear dynamics of learning in deep linear neural networks. In Yoshua Bengio and Yann LeCun (eds.), 2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014, Conference Track Proceedings, 2014. URL http://arxiv.org/abs/1312.6120.
- Andrew M Saxe, James L McClelland, and Surya Ganguli. A mathematical theory of semantic development in deep neural networks. *Proceedings of the National Academy of Sciences*, 116(23): 11537–11546, 2019.
- James B Simon, Madeline Dickens, Dhruva Karkada, and Michael R DeWeese. The eigenlearning framework: A conservation law perspective on kernel regression and wide neural networks. *Transactions on Machine Learning Research*, 2023. ISSN 2835-8856.
- Stefano Spigler, Mario Geiger, and Matthieu Wyart. Asymptotic learning curves of kernel methods: empirical data versus teacher–student paradigm. *Journal of Statistical Mechanics: Theory and Experiment*, 2020(12):124001, 2020.
- James Stevens et al. *Applied multivariate statistics for the social sciences*, volume 4. Lawrence erlbaum associates Mahwah, NJ, 2002.
- Jingtong Su, Ya Shi Zhang, Nikolaos Tsilivis, and Julia Kempe. On the robustness of neural collapse and the neural collapse of robustness. *arXiv preprint arXiv:2311.07444*, 2023.

- Peter Súkeník, Marco Mondelli, and Christoph Lampert. Neural collapse versus low-rank bias: Is deep neural collapse really optimal? *arXiv preprint arXiv:2405.14468*, 2024.
- Salma Tarmoun, Guilherme Franca, Benjamin D Haeffele, and Rene Vidal. Understanding the dynamics of gradient flow in overparameterized linear models. In *International Conference on Machine Learning*, pp. 10153–10161. PMLR, 2021.
- Michael te Vrugt and Raphael Wittkowski. Projection operators in statistical mechanics: a pedagogical approach. *European Journal of Physics*, 41(4):045101, 2020.
- Christopher KI Williams and Carl Edward Rasmussen. *Gaussian processes for machine learning*, volume 2. MIT press Cambridge, MA, 2006.
- Jing Xu and Haoxiong Liu. Quantifying the variability collapse of neural networks. In *International Conference on Machine Learning*, pp. 38535–38550. PMLR, 2023.
- Greg Yang and Edward J. Hu. Tensor programs IV: Feature learning in infinite-width neural networks. In *Proceedings of the 38th International Conference on Machine Learning*, pp. 11727–11737. PMLR, 2021. URL https://proceedings.mlr.press/v139/yang21c.html. ISSN: 2640-3498.
- Zhihui Zhu, Tianyu Ding, Jinxin Zhou, Xiao Li, Chong You, Jeremias Sulam, and Qing Qu. A geometric analysis of neural collapse with unconstrained features. *Advances in Neural Information Processing Systems*, 34:29820–29834, 2021.

# A GLOSSARY

756

Symbol.	Name	Definition	Ref
$\overline{C}$	Class count	The number of classes.	Section 2
p	Layer width	The width of the last layer	Section 2
n	Sample count	The number of training samples	Section 2
$\mathcal{X}$	Input space	Space of inputs	Section 2
q	Input distribution	The probability distribution that generates samples in the input space.	Section 2
Φ	Feature map	A map from input to post-activation of the penultimate layer (the activations fed to the last layer).	Section 2
$\Phi(x)$	(Last layer) Features	A $p$ -dimensional random variable or the post-activation of penultimate layer for $x\sim q$ .	Section 2
$\mathcal{T}$	(Feature kernel) integral operator	A map from Hilbert space to Hilbert space that depends on $\Phi(x)$ , $x \sim q$ (typically the last layer features)	Eq. (1)
$\mathcal{T}_{MP}$	Minimum projection operator	A special set of operators that depends on the learned function $\hat{f}$ .	Definition 1
$e_k$	$k^{\rm th}$ eigenfunction	The eigenfunction of $\mathcal{T}$ .	Eq. (2)
$ ho_k$	$k^{\mathrm{th}}$ eigenvalue	The eigenvalue of $\mathcal{T}$ .	Eq. (2)
$\hat{f}$	Learned function	A $C$ -dimensional vector output function $\hat{f}: \mathcal{X} \to \mathbb{R}^C$ expressed by the neural network.	Section 2
$f^*$	Target function	A $C$ -dimensional vector output function $f^*: \mathcal{X} \to \mathbb{R}^C$ with correct labels. The output is always a one-hot vector (up to a scaling constant).	Section 2
$\mathcal{H}^*$	Target function space	Space linearly spanned by the entries of the target function.	Section 3
$\hat{\mathcal{H}}$	Learned function space	Space linearly spanned by the entries of the learned function	Section 3
$P_{\mathcal{H}}$	Projection operator	Projection operator onto $\mathcal{H}$ . If $\mathcal{H} = span\{e_1, e_2, \dots, e_p\}$ , where $e_k$ 's are orthonormal, the the projection operator is given as $P_{\mathcal{H}} = \sum_{k=1}^{p}  e_k\rangle \langle e_k $ .	Section 3
$\alpha$	Downscale constant	Prefactor that downscaling the target such that $y \to y/\alpha$	Section 4.1
$\Sigma_b$	Inter-class covariance matrix	Covariance of class mean vectors for the training set.	Eq. (33)
$\Sigma_W$	Intra-class covariance matrix	Covariance of training feature vectors within given class.	Eq. (32)
$CKA(\cdot, \cdot)$	Centered kernel alignment	Alignment measure between two matrices or operators. It is normalized between $[0,1]$ and ignores the mean (for matrices) or constant function (for operators)	Eq. (6)
$\ \theta\ _F^2$	Parameter norm	Norm of all parameters in the model. It was used as richness measure with smaller value meaning richer dynamics.	Section 4.1
$\mathcal{D}_{LR}$	low rank metric	Our proposed metric of dynamical richness. Smaller is richer.	Eq. (4)
$\mathcal{S}_{init}$	Kernel deviation	The CKA measure between the learned and initial kernel.	Section 4.1
NC1	Neural collapse metric	The trace of $\Sigma_b^{\dagger}\Sigma_W$ , measuring the training feature vectors variance compared to class boundaries.	Proposition 1
$\Pi^*(k)$	Cumulative quality	Measure of how well the first $k$ eigenfunctions span the target function space.	Eq. (5)
$\hat{\Pi}(k)$	Cumulative utilization	Measure of how well the first $k$ eigenfunctions span the learned function space.	Eq. (5)
$Q^*(k)$	Quality	Per-feature quality or $\Pi^*(k) - \Pi^*(k-1)$ where $\Pi^*(0) = 0$	Fig. 5
$\hat{Q}(\hat{k})$	Utilization	Per-feature utilization or $\hat{\Pi}(k) - \hat{\Pi}(k-1)$ where $\hat{\Pi}(0) = 0$	Fig. 5

# B TECHNICAL SUPPLEMENTARY MATERIAL

Here, we introduce technical terms used in the main text and the following appendices.

#### B.1 Bra-ket notation

In physics, bra-ket notations are widely used to express the inner product in function space. In our paper, we use them to avoid overload of expectations and to clarify that we are using functions. Note that we use bra-ket notation for the expectation over the true input distribution only:

$$\langle f|g\rangle = \mathbf{E}_{x\sim g}[f(x)g(x)].$$

The notation is also useful for expressing operators such as  $\mathcal{T}$ :

$$\mathcal{T}[f](x) = \sum_{k=1}^{p} |\Phi_k\rangle \langle \Phi_k|f\rangle = [\Phi_1 \langle \Phi_1|f\rangle, \dots, \Phi_p \langle \Phi_p|f\rangle].$$

Because  $\mathcal{T}$  maps vector function  $f: \mathcal{X} \to \mathbb{R}^p$  to a vector function,  $\mathcal{T}[f]: \mathcal{X} \to \mathbb{R}^p$  is also a vector function with

$$\mathcal{T}[f](x) = [\Phi_1(x) \langle \Phi_1 | f \rangle, \dots, \Phi_p(x) \langle \Phi_p | f \rangle],$$

where  $\langle \Phi_k | f \rangle$  are scalars and  $\Phi_k : \mathcal{X} \to \mathbb{R}$  are functions.

# B.2 FEATURES ARE RANDOM VARIABLES AND (WELL-BEHAVED) RANDOM VARIABLES FORM A HILBERT SPACE

In machine learning, textbooks often overlook the mathematical distinction between feature maps and features, treating them as interchangeable. However, they are fundamentally different.

A feature map is a function  $f: \mathcal{X} \to \mathcal{Y}$ , defined independently of any distribution. A feature, by contrast, is a random variable induced by applying a feature map to inputs drawn from a distribution q over  $\mathcal{X}$ .

Features are therefore distribution-dependent: applying the same feature map to different input distributions yields different features. For instance, a fixed neural network defines a feature map, but the resulting features — such as last-layer activations — will differ between MNIST and Fashion-MNIST due to changes in the input distribution. The distinction also applies to kernel operators  $\mathcal{T}:\mathcal{H}\to\mathcal{H}$  (distribution dependent) and kernels  $K:\mathcal{X}\times\mathcal{X}\to\mathbb{R}$  (distribution independent).

The distinction becomes more important when we wish to discuss Hilbert space. A Hilbert space requires an inner product between functions (feature maps), which depends on the underlying distribution. For example,  $\sin(x)$  and  $\sin(2x)$  are orthogonal for  $q = unif[0, 2\pi]$  where  $\int_0^{2\pi} \sin(x) \sin(2x) dx = 0$ , but not for unit Gaussian distribution  $q = \mathcal{N}(0, 1)$  where  $\int_{-\infty}^{\infty} \sin(x) \sin(2x) e^{-x^2/2} dx \neq 0$ .

For an underlying probability distribution and a set of functions (feature maps), we can form a Hilbert space. A common example of a Hilbert space will be the set of solutions expressed by a linear model such as

$$f(x) = w_1 x + w_2 x^2 + w_3 x^3$$
.

Assuming the features  $[x, x^2, x^3]$  are linearly independent on input distribution q, they form a 3-dimensional Hilbert space.

# B.3 CENTERED KERNEL ALIGNMENT

The centered kernel alignment (CKA) (Kornblith et al., 2019) measures the similarity between two matrices or operators, and are commonly used for tracking the evolution of NTK (Baratin et al., 2021; Lou et al., 2022):

$$CKA(A,B) = \frac{\text{Tr}(c(A)c(B))}{\|c(A)\|_F \|c(B)\|_F},$$
(6)

where  $c(A) = (I - |\mathbf{1}\rangle \langle \mathbf{1}|) A (I - |\mathbf{1}\rangle \langle \mathbf{1}|)$  is centering operator, where I is identity,  $|\mathbf{1}\rangle$  is the constant function or constant vector, and  $\|\cdot\|_F$  is Frobenius norm. The centering operator removes the constant shift. The analogy is measuring the similarity of two multivariate Gaussian random variables by comparing their covariance matrices  $\mathbf{E}[(X - \mathbf{E}[X])(X - \mathbf{E}[X])^T]$  instead of their autocorrelation matrices  $\mathbf{E}[XX^T]$ .

For  $\mathcal{T}_{MP}$ , if the learned function space  $\hat{\mathcal{H}}$  contains the constant function — which is empirically true throughout the training for all our experiments and always true if it perfectly fits the training samples (Appendix D) —  $c(\mathcal{T}_{MP})/\|c(\mathcal{T}_{MP})\|_F$  becomes an orthogonal projection operator on the constant complement of  $\hat{\mathcal{H}}$ , and is independent of the values of  $a_1$  and  $a_2$ .

#### B.4 EFFECTIVE DIMENSION

A covariance matrix or an integral operator may be fully ranked, but their eigenvalues may decay fast (e.g., by a power-law). The small eigendirections are often negligible, and researchers use various effective dimensions to measure the number of **significant** dimensions. In this appendix, we will use the exponent of entropy (Hill, 1973) to measure the effective dimension.

For a matrix or an operator with positive eigenvalues  $[\rho_1, \cdots, \rho_p]$ , the effective dimension is

$$\exp\left(-\sum_{i=1}^{p} \frac{\rho_i}{\sum_{j=1}^{p} \rho_j} \ln\left(\frac{\rho_i}{\sum_{j=1}^{p} \rho_j}\right)\right). \tag{7}$$

Note that the effective dimension, the exponential of Shannon entropy, is d when the eigenvalues have d equal non-zero entries. A slower decay of entries of  $\rho$  (in non-increasing order of entries) generally yields a higher effective dimension than a vector with a faster decay.

The effective dimension of a matrix (operator) can be interpreted as the number of linearly independent vectors (functions) needed to effectively describe the matrix (operator). This is similar in spirit to the principal component analysis (PCA) (Abdi & Williams, 2010) in that only the directions with significant variance are considered.

As discussed in Appendix D, the first eigenfunction is always a constant function, which is irrelevant for our measure  $\mathcal{D}_{LR}$  and dynamical interpretation. Because the first eigenvalue is often much larger than other significant eigenvalues, we use

$$D_{eff}(\rho) := 1 + \exp\left(-\sum_{i=2}^{p} \frac{\rho_i}{\sum_{j=2}^{p} \rho_j} \ln\left(\frac{\rho_i}{\sum_{j=2}^{p} \rho_j}\right)\right). \tag{8}$$

#### C LINEAR MODELS, FEATURES, AND THEIR INDUCTIVE BIAS

Here, we will introduce the recent findings on the inductive bias of linear models and demonstrate the significance of eigenvalues and eigenfunctions of  $\mathcal{T}$ . Linear regression or kernel regression is a rare case where its inductive bias is analytically calculable (Bordelon et al., 2020; Jacot et al., 2020; Spigler et al., 2020; Canatar et al., 2021; Cui et al., 2021; El Harzli et al., 2024; Simon et al., 2023). For a given feature map  $\Phi: \mathcal{X} \to \mathbb{R}^p$  or kernel  $K(x, x') = \Phi(x)^T \Phi(x')$ , the model is expressed as the p dimensional (p can be infitnite) linear model

$$\hat{f}(x;w) = \sum_{k=1}^{p} w_k \Phi_k(x). \tag{9}$$

The objective of (ridgeless) regression is to minimize the following empirical loss

$$\mathcal{L}_{emp}(\hat{f}) = \frac{1}{2} \sum_{i=1}^{n} \left| \hat{f}(x^{(i)}) - f^*(x^{(i)}) \right|^2 + \lim_{\lambda \to 0} \lambda ||w||_F^2, \tag{10}$$

where  $x^{(i)}$  denotes the  $i^{\rm th}$  sample in the training set. The second term is the regularization which we take the limit so the solution is unique (pseudo inverse solution) for an overparameterized setup (n < p).

#### C.1 EIGENFUNCTIONS - ORTHONORMALIZED FEATURES

In many cases,  $[\Phi_1(x), \Phi_2(x), \cdots, \Phi_p(x)]$  are not orthonormal. For example, the features  $[x, x^2]$  of  $f(x) = w_1 x + w_2 x^2$  for q = unif[0, 1] are not orthonormal:

$$\langle x|x^2\rangle = \int_0^1 x^3 dx \neq 0, \quad \langle x|x\rangle = \int_0^1 x^2 dx \neq 1, \quad \langle x^2|x^2\rangle = \int_0^1 x^4 dx \neq 1. \tag{11}$$

The diagonalization into the eigenfunctions returns the orthonormal basis of the hypothesis (Hilbert) space  $\mathcal{H}=span\{\Phi_1(x),\Phi_2(x),\cdots,\Phi_p(x)\}$ , more formally known as Reproducing Kernel Hilbert Space (RKHS). The eigenvalue now shows the norm of the feature along the direction of eigenfunction as the features  $[\Phi_1(x),\Phi_2(x),\cdots,\Phi_p(x)]$  are not normalized (Fig. 6).

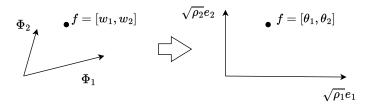


Figure 6: **Diagonalization of features** Since the features  $[\Phi_1(x), \Phi_2(x), \cdots, \Phi_p(x)]$  span a vector space, each function  $\Phi_k : \mathcal{X} \to \mathbb{R}$  can be represented as a vector. The overlap among  $\Phi_k$ 's (non-zero linear correlation) and differing norms of  $\Phi_k$ 's result in diagonalized features  $[\sqrt{\rho_1}e_1, \sqrt{\rho_2}e_2, \cdots, \sqrt{\rho_p}e_p]$  to have distinct intensities (norms)  $\rho_1, \cdots, \rho_p$ .

The transformation between the eigenfunction basis and the feature basis is

$$\Phi_i(x) = \sum_j O_{ij} \sqrt{\rho_j} e_j(x), \qquad O_{ij} = \frac{1}{\sqrt{\rho_j}} \langle \Phi_i | e_j \rangle, \qquad (12)$$

where  $O \in \mathbb{R}^{p \times p}$  is an orthogonal matrix (follows trivially from Eq. (2)). Using Eq. (12), the model can be reparameterized in the eigenfunction basis

$$f(x) = \sum_{k=1}^{p} w_k \Phi_k(x) = \sum_{k=1}^{p} \theta_k \sqrt{\rho_k} e_k(x).$$
 (13)

Note that because  $\theta = Ow$ , the norm of the parameters is conserved (i.e.  $\|\theta\|_F = \|w\|_F$ ), indicating that weighted basis  $[\sqrt{\rho_1}e_1, \sqrt{\rho_2}e_2, \cdots, \sqrt{\rho_p}e_p]$  instead of normalized basis  $[e_1, e_2, \cdots, e_p]$  must be used to represent the 'intensity' or norm of the features  $[\Phi_1(x), \cdots, \Phi_p(x)]$ .

#### C.2 INDUCTIVE BIAS TOWARD LARGE FEATURES FOR MINIMUM NORM SOLUTION

For overparameterized linear models, the inductive bias determines the returned function as many functions can express the training set. Here, we brief the linear models' inductive bias toward larger (significant) features or eigenfunctions (Bordelon et al., 2020; Jacot et al., 2020; Spigler et al., 2020; Canatar et al., 2021; Cui et al., 2021; El Harzli et al., 2024; Simon et al., 2023).

The inductive bias is quantified through learnability  $L_k$  (Simon et al., 2023) or the expectation ratio between the learned coefficient and true coefficient for  $e_k$ :

$$L_k := \mathbf{E}_{S \sim q^n} \left[ \frac{\langle \hat{f} | e_k \rangle}{\langle f^* | e_k \rangle} \right] = \frac{\rho_k}{\rho_k + \kappa}, \quad \text{where } \sum_{k=1}^p \frac{\rho_k}{\rho_k + \kappa} = n.$$
 (14)

The learnabilities always sum to n — the number of training datapoints — and constant  $\kappa$  is the constant that satisfies the equality.

We can understand the inductive bias from the (vanishing) regularization on  $L^2$  norm  $(\lambda ||w||_F^2$  in Eq. (10)) and parameterization in Eq. (13). If different eigenfunctions can equally express the training set, eigenfunction with larger eigenvalue often requires smaller coefficient  $\theta_k$  to express the samples. Because  $||w||_F^2 = ||\theta||_F^2$ , expressing the data with larger (eigenvalue) eigenfunctions minimizes the norm, creating an inductive bias toward larger features.

Instead of the derivation found in the references (Bordelon et al., 2020; Jacot et al., 2020; Spigler et al., 2020; Canatar et al., 2021; Cui et al., 2021; El Harzli et al., 2024) — which requires random matrix theory or replica trick — we show an example of the inductive bias in Fig. 7.

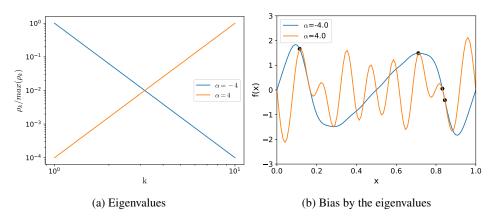


Figure 7: **Different inductive biases by the eigenvalues (intensity).** Two 10-parameter linear models are trained on 4 datapoints (overparmaterized) with gradient flow. Both linear models use  $\sin$  basis functions such that  $f(x) = \sum_{k=0}^{10} w_k k^{\alpha} \sqrt{2} \sin(2\pi kx)$  — spanning an identical function space — but differ in the eigenvalues with  $\alpha = -4$  (blue) and  $\alpha = 4$  (orange), leading to different inductive biases. (a): The blue model has a greater intensity for lower-frequency  $\sin$  functions, while the orange model shows the opposite. (b): The learned functions show that the blue model used lower-frequency functions while the orange used higher-frequency functions to express the training samples.

Fig. 7 shows that two linear models with same hypothesis space but different eigenvalues learn dramatically different functions. The blue model has large eigenvalues for low-frequency functions and fits the datapoints using low-frequency functions. The orange model with opposite eigenvalue distribution fits the training set with high-frequency functions — a clear example of an inductive bias toward large (intensity) features.

In the references, the generalization loss is

$$\mathcal{L}_{G} = \frac{1}{1 - \frac{1}{n} \sum_{j}^{p} L_{j}^{2}} \left( \sum_{k}^{p} (1 - L_{k})^{2} \left\langle e_{k} | f^{*} \right\rangle^{*2} \right). \tag{15}$$

The equation formalizes the task-model alignment: the alignment between large features (large  $\rho_k$  thus large  $L_k$ ) and high quality (large  $\langle e_k | f^* \rangle$ ) leads to better generalization.

#### C.3 DYNAMICAL INDUCTIVE BIAS TOWARD LARGE FEATURES

Here, we review how gradient descent introduces the dynamical inductive bias toward larger features. For a linear model with feature map  $\Phi$  and kernel integral operator  $\mathcal{T}$ , the gradient flow dynamics under MSE loss with infinite training datapoints and zero initialization becomes

$$\frac{df}{dt} = -\mathcal{T}[f - f^*]. \tag{16}$$

The derivation is trivial with

$$\frac{df}{dt}(x') = \sum_{k=0}^{p} \frac{dw_k}{dt} \Phi_k(x') = \sum_{k=0}^{p} -\frac{\partial \mathcal{L}}{\partial w_k} \Phi_k(x')$$
(17)

$$= -\mathbf{E}_x \left[ (f(x) - f^*(x)) \sum_{k=0}^{p} \frac{df}{dw_k}(x) \right] \Phi_k(x')$$
(18)

$$= -\mathbf{E}_x \left[ (f(x) - f^*(x)) \sum_{k=0}^{p} \Phi_k(x) \right] \Phi_k(x')$$
(19)

$$= -\mathcal{T}\left[f - f^*\right],\tag{20}$$

where in the second line, we used Eq. (10) on infinite datapoints so  $\mathcal{L}_{emp} = \mathcal{L}$ . In the last line, we used the definition of integral operator (Eq. (1)).

Using the eigenvalues  $\rho_k$  and eigenfunctions  $e_k$  of the kernel integral operator  $\mathcal{T}$ , we obtain the dynamics for each  $\langle e_k | f \rangle$  by inner producting  $e_k$  on both sides of Eq. (16),

$$\frac{d\langle e_k|f\rangle}{dt} = \mathbf{E}_{x'} \left[ e_k(x') \frac{df}{dt}(x') \right] = -\mathbf{E}_{x'} \left[ e_k(x') \mathcal{T} \left[ f - f^* \right] (x') \right]$$
(21)

$$= -\langle e_k | \mathcal{T}[f - f^*] \rangle \tag{22}$$

$$= -\rho_k(\langle e_k | f \rangle - \langle e_k | f^* \rangle), \tag{23}$$

where we use the definition of eigenfunction in the last line (Eq. (2)). Expanding f in the eigenfunction basis and plugging in Eq. (21), the function f at time t is a sum of p independent modes that saturate faster for larger eigenvalues:

$$f(x)|_{t} = \sum_{k=1}^{p} \langle e_{k}|f^{*}\rangle (1 - e^{-\rho_{k}t})e_{k}(x),$$
 (24)

where we assumed f is a zero function at initialization. Eq. (24) shows that gradient flow decouples the dynamics of linear models into p modes, where each mode corresponds to the evolution of  $\langle e_k|f\rangle$  having a saturation speed governed by  $\rho_k$  — thus the dynamical inductive bias toward larger eigenfunction.

#### C.4 APPLICATION TO OUR METHOD

Here, we detail the intuition of our visualization methods. As described in Bengio et al. (2013), feature learning (of representation) is providing **better** features for a simpler learner. The analogy allows us to describe a neural network as (1) feature map  $\Phi$  providing better features for (2) the last layer — the simple learner (Fig. 8). Yet, **better** features for the last layer is ill-defined.

The last layer interacts with rest of the network only through the features. While the exact dynamics differs from linear models, the last layer maintains the dynamical inductive bias (Appendix C.3) toward larger features. Furthermore, the quality  $(\langle e_k|f^*\rangle)$ , utilization  $(\langle e_k|\hat{f}\rangle)$ , and intensity  $\rho_k$  formalisms readily extend to the last layer features, motivating our visualization method.

In the rich regime, features evolve during training, and Fig. 9(a-c) demonstrates an example of change in features for 4-layer MLP trained to fit the Heaviside step function. The eigenfunctions are difficult to visualize in practice as the input space is high dimensional. Fig. 9(d-f) shows our visualization, capturing the evolution of features by measuring properties natural to linear regression.

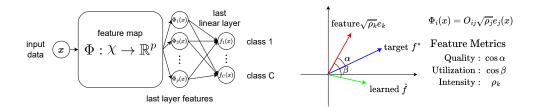


Figure 8: A neural network decomposed into a feature map and a linear last layer. (left): An abstract diagram depicting a DNN architecture as a combination of the feature map  $\Phi: \mathcal{X} \to \mathbb{R}^p$  from the input space  $\mathcal{X}$  to the p post activations of the penultimate layer, and final linear classifier for C-way classification. Most neural networks share this abstract structure and mainly vary in how they create their feature maps. (right): Illustration of the visualization methods in function space.

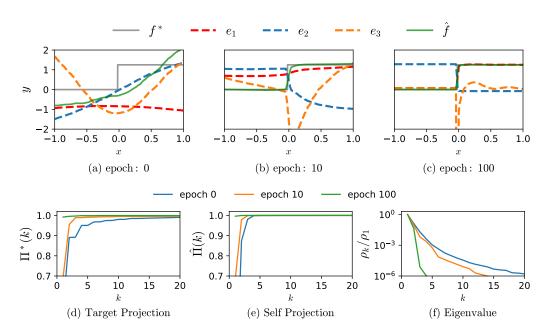


Figure 9: **Toy model demonstrating feature learning.** A 4-layer MLP with width 1000 and scalar input and output is trained to learn the Heaviside step function  $f^*$  over the domain [-1,1]. (a): The first three eigenfunctions (dashed) are shown at initialization, resembling orthogonal polynomials. (b,c): During training, features evolve such that a single feature (red) fits both the target and learned functions (grey and green). (d,e,f): Our visualization metrics defined in Eq. (5) at different epochs. (d) and (e) show that fewer features can express the target/learned function, and (f) shows that only a few features are significant after training.

# D THE RELATIONSHIP BETWEEN THE MP-OPERATOR $\mathcal{T}_{\mathbf{MP}}$ AND NEURAL COLLAPSE

Neural collapse (Papyan et al., 2020) (NC) refers to a state of a DNN when the feature vectors of the training set and last layer weights form a symmetric and clustered structure at the terminal phase of training (TPT) or when trained past the point where the training error vanishes. The structure has been studied mainly in the Unconstrained Feature Model (UFM) (Mixon et al., 2020) and has sparked theoretical studies. See Kothapalli et al. (2022) for a review.

NC is defined by the emergence of four interconnected phenomena upon TPT: NC1) collapse of within-class variability, NC2) convergence of features to a rigid simplex equiangular tight frame (ETF) structure, NC3) alignment of the last layer and features, and NC4) simplified decision by nearest class:

#### 1. Within class variance tends to 0

$$\Sigma_W \Sigma_h^{\dagger} \to 0,$$
 (25)

### 2. Convergence to simplex ETF

$$\frac{(\mu_i - \bar{\mu})^T (\mu_j - \bar{\mu})}{\|(\mu_i - \bar{\mu})\|_2 \|(\mu_j - \bar{\mu})\|_2} \to \frac{C\delta_{ij} - 1}{C - 1},\tag{26}$$

#### 3. Convergence to self duality

$$\frac{w_i}{\|w_i\|_2} - \frac{\mu_i - \bar{\mu}}{\|(\mu_i - \bar{\mu})\|_2} \to 0, \tag{27}$$

#### 4. Simplification to nearest class center

$$\arg\max_{i} w_i \Phi(x) + b_i \to \arg\min_{i} \|\Phi(x) - \mu_i\|_2.$$
 (28)

#### D.1 Proofs of Proposition 1 and 2

For completeness, we restate the conditions. We assume the true distribution q equals the empirical distribution of the training set — the expectation  $\mathbf{E}_x$  (and bra-ket notation) is over the training samples. We assume balanced classification and define  $A_i$  as the set of training samples with class label i. We assume the learned function is a perfectly classifying indicator function, giving

$$\hat{f}_i(x) := \begin{cases} 0 & : x \notin A_i \\ 1 & : x \in A_i \end{cases}$$
 (29)

The feature class mean vector for label i is

$$\mu_i = \mathbf{E}_{x \in A_i}[\Phi(x)] = \frac{C}{n} \sum_{x \in A_i} \Phi(x) = C \langle \hat{f}_i | \Phi \rangle.$$
 (30)

Note that  $\langle \hat{f}_i | \Phi \rangle \in \mathbb{R}^p$  is a vector with  $\langle \hat{f}_i | \Phi \rangle = [\langle \hat{f}_i | \Phi_1 \rangle, \dots, \langle \hat{f}_i | \Phi_p \rangle].$ 

The global mean vector is

$$\bar{\mu} = \frac{1}{C} \sum_{i}^{C} \mu_{i} = \langle \Phi | \sum_{i=1}^{C} \hat{f}_{i} \rangle = \langle \mathbf{1} | \Phi \rangle.$$
 (31)

The feature intra (within)-class covariance  $\Sigma_W \in \mathbb{R}^{p \times p}$  is

$$\Sigma_W := \sum_{i=1}^{C} \mathbf{E}_{x \in A_i} \left[ (\Phi(x) - \mu_i)(\Phi(x) - \mu_i)^T \right]. \tag{32}$$

The feature inter (between)-class covariance  $\Sigma_b \in \mathbb{R}^{p \times p}$  is

$$\Sigma_b := \frac{1}{C} \sum_{i=1}^{C} (\mu_i - \bar{\mu})(\mu_i - \bar{\mu})^T.$$
 (33)

**Lemma 1.** For the conditions given above, if  $\mathcal{T}[f](x') = \Phi(x')^T \mathbf{E}_{x \sim q'}[\Phi(x)f(x)]$  is  $\mathcal{T}_{MP}$ , then

$$\Phi^{T}(x')^{T}(\mu_{i} - \bar{\mu}) = a_{2}(C\hat{f}_{i}(x') - 1). \tag{34}$$

*Proof.* By the definition of indicator functions of equal partitions,

$$\langle \hat{f}_i | \hat{f}_j \rangle = \frac{1}{C} \delta_{ij}, \qquad \sum_{i=1}^C \hat{f}_i = \mathbf{1}, \qquad \langle \mathbf{1} | \hat{f}_j \rangle = \frac{1}{C}.$$
 (35)

It follows that  $C\hat{f}_j - \mathbf{1}$  is orthogonal to 1:

$$\langle \mathbf{1}|C\hat{f}_i - \mathbf{1}\rangle = 0. \tag{36}$$

We can express the function  $\Phi(x')^T(\mu_i - \bar{\mu}) : \mathcal{X} \to \mathbb{R}$  in terms of  $\mathcal{T}$ 

$$\Phi(x')^{T}(\mu_{i} - \bar{\mu}) = \Phi(x')^{T} \left( \langle \Phi | C\hat{f}_{i} \rangle - \langle \Phi | \mathbf{1} \rangle \right)$$
(37)

$$= \Phi(x')^T \mathbf{E}_x [\Phi(x) (C\hat{f}_i(x) - 1)]$$
 (38)

$$= \mathcal{T}[C\hat{f}_i - 1](x'), \tag{39}$$

where we used the definition of the means (Eqs. (30) and (31)) in the first line and the definition of  $\mathcal{T}$  (Eq. (1)) in the last line. As  $\mathcal{T}$  is  $\mathcal{T}_{MP}$ , we obtain

$$\Phi(x')^{T}(\mu_{i} - \bar{\mu}) = \mathcal{T}_{MP}[C\hat{f}_{i} - 1](x')$$
(40)

$$= a_2(C\hat{f}_i(x') - 1), \tag{41}$$

where we used Eq. (36) and the definition of  $\mathcal{T}_{MP}$  (Definition 1).

**Corollary 1.** *In the setup of Lemma 1, we have* 

$$\mu_i^T(\mu_j - \bar{\mu}) = a_2(C\delta_{ij} - 1).$$

*Proof.* From the definition of class mean (Eq. (30)) and Lemma 1,

$$\mu_j^T(\mu_i - \bar{\mu}) = \mathbf{E}_{x'} \left[ C\hat{f}_j(x')\Phi(x')^T(\mu_i - \bar{\mu}) \right]$$
(42)

$$= \mathbf{E}_{x'} \left[ C\hat{f}_j(x') a_2 (C\hat{f}_i(x') - 1) \right]$$

$$\tag{43}$$

$$= a_2(C\delta_{ij} - 1), \tag{44}$$

where we used that  $\mathbf{E}_{x'}[\hat{f}_{i}^{2}(x')] = \mathbf{E}_{x'}[\hat{f}_{i}(x')] = C^{-1}$ .

**Proposition 1.** If  $\mathcal{T}$  is an MP-operator, then the NC1 condition (collapse of within-class variability)  $\Sigma_b^{\dagger} \Sigma_W = 0$  holds, where inter-class covariance matrix  $\Sigma_b$  and intra-class covariance matrix  $\Sigma_W$  are

$$\Sigma_b = \sum_{j=1}^C (\mu_j - \bar{\mu})(\mu_j - \bar{\mu})^T \quad \text{and} \quad \Sigma_W = \sum_{i=1}^C \mathbf{E}_{x \in A_i} \left[ (\Phi(x) - \mu_i)(\Phi(x) - \mu_i)^T \right].$$

*Proof.* We will show that  $Tr(\Sigma_b\Sigma_W)=0$ , which automatically proves the proposition as covariance matrices are positive semi-definite. We have

$$Tr(\Sigma_b \Sigma_W) = Tr\left(\sum_j (\mu_j - \bar{\mu})(\mu_j - \bar{\mu})^T \sum_i \mathbf{E}_{x \in A_i} (\Phi(x) - \mu_i)(\Phi(x) - \mu_i)^T\right)$$
(45)

$$= \sum_{i,j} \mathbf{E}_{x \in A_i} \left[ \left( (\Phi(x) - \mu_i)^T (\mu_j - \bar{\mu}) \right)^2 \right]$$

$$\tag{46}$$

$$= \sum_{i,j} \mathbf{E}_{x \in A_i} \left[ a_2^2 \left( C \hat{f}_j(x) - 1 - C \delta_{ij} + 1 \right)^2 \right]$$
 (47)

$$= a_2^2 C^2 \sum_{i,j} \mathbf{E}_{x \in A_i} \left[ \left( \hat{f}_j(x) - \delta_{ij} \right)^2 \right], \tag{48}$$

where we used linearity of trace, sum, expectation in the first line and Lemma 1 in the second line. Expanding Eq. (48),

$$Tr(\Sigma_b \Sigma_W) = a_2^2 C^2 \sum_{i,j} \mathbf{E}_{x \in A_i} \left[ (1 - 2\delta_{ij}) \hat{f}_j(x) + \delta_{ij} \right]$$
(49)

$$= a_2^2 C \sum_{i,j} (\delta_{ij} (2 - 2\delta_{ij})) \tag{50}$$

$$=0, (51)$$

where we used that learned functions are indicator functions in the second line.

**Proposition 2.** If  $\mathcal{T}$  is an MP-operator, then the NC2 condition (convergence of features to a simplex equiangular tight frame) holds:

$$(\mu_i - \bar{\mu})^T (\mu_j - \bar{\mu}) \propto \delta_{ij} - \frac{1}{C}.$$

*Proof.* From the definition of global mean (Eq. (31)) and Lemma 1,

$$\bar{\mu}^T(\mu_i - \bar{\mu}) = \mathbf{E}_{x'} \left[ \Phi(x')^T (\mu_i - \bar{\mu}) \right] \tag{52}$$

$$= \mathbf{E}_{x'} \left[ a_2 (C\hat{f}_i(x') - 1) \right] \tag{53}$$

$$=0. (54)$$

Using Corollary 1 and Eq. (54),

$$(\mu_i - \bar{\mu})^T (\mu_j - \bar{\mu}) = Ca_2(\delta_{ij} - \frac{1}{C}).$$
 (55)

## THE FIRST EIGENFUNCTION, THE CONSTANT FUNCTION, THE ANALOG OF GLOBAL MEAN VECTOR IN SIMPLEXT ETF

Here, we discuss the constant function, which, by definition of  $\mathcal{T}_{MP}$  (Definition 1), is the largest eigenfunction if  $\mathcal{T}$  is  $\mathcal{T}_{MP}$ . In all experiments except Fig. 9, the largest eigenfunction is the constant function throughout training, with  $\langle \mathbf{1} | e_1 \rangle > 0.95$ . For our visualization, we ignore the first eigenvalue from analysis as the constant function is removed in CKA (Eq. (6)).

In the NC2 condition, the k-simplex ETF structure is a set of k orthogonal vectors projected on k-1dimensional space along the compliment of the global mean vector. For example, projecting [0,0,1], [0,1,0], and [1,0,0] onto the orthogonal complement of  $[\frac{1}{3},\frac{1}{3},\frac{1}{3}]$  (the global mean vector) gives three vertices  $\left[-\frac{1}{3}, -\frac{1}{3}, \frac{2}{3}\right]$ ,  $\left[-\frac{1}{3}, \frac{2}{3}, -\frac{1}{3}\right]$ ,  $\left[\frac{2}{3}, -\frac{1}{3}, -\frac{2}{3}\right]$  which forms an (ordinary plane) equilateral triangle.

Our operator  $\mathcal{T}_{MP}$  (Definition 1) has an analogous structure, where it is a (scaled) projection operator  $(a_2P_{\hat{H}})$  except along the constant function. The CKA (Eq. (6)) also measures the alignment except along the mean direction (for matrices) or the constant direction (for operators). This allows neural networks to have any arbitrary eigenvalue for the constant function, yet perfectly align with  $\mathcal{D}_{LR}=0$ .

When  $\mathcal{T}$  is  $\mathcal{T}_{MP}$ , the definition of  $\mathcal{T}_{MP}$  (Definition 1) trivially shows that the largest eigenfunction is the constant function. A peculiar observation is that the first eigenfunction of  $\mathcal T$  remains the constant function ( $\langle 1|e_1\rangle > 0.95$ ) throughout training in all our experiments (except Fig. 9). This pattern also appears in hierarchical datasets (Saxe et al., 2022), suggesting that natural data may consistently prioritize the constant function as the dominant mode — warranting further investigation.

#### GENERALITY OF OUR METRIC BEYOND NEURAL COLLAPSE

As discussed in the main text, our metrics are defined in terms of functions (random variables) while neural collapse is defined by the feature vectors of the training set. Working in the function space, we can measure the quality Eq. (5), which is inaccessible from formalism using training feature vectors.

 Another important difference is that we focus on the learned function instead of mean class label vectors. This allows more independence beyond neural collapse, allowing the method to be applied on broader set of tasks: tasks with orthogonal and isotropic target functions form a broader class than balanced classification tasks. Regression task with scalar output is an example.

In Fig. 10, we train ResNet18 on MNIST as a **regression** task with scalar output  $(f^*: \mathcal{X} \to \mathbb{R})$ , where a digit i is correctly predicted if i-0.5 < f(x) < i+0.5. Unlike C-way classification, which typically uses C features, the regression model reaches the rich regime with just two — one of which is constant. Neural collapse formalism, focusing on class labels, does not readily extend to regression tasks, demonstrating that our function-based formalism is more general, label-independent, and applicable beyond classification.

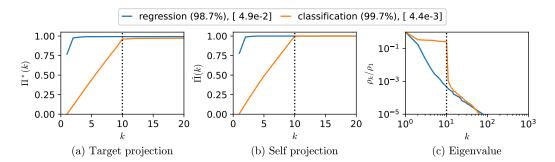


Figure 10: Rich dynamics for regression problem with scalar output. ResNet18 was trained on MNIST via regression (blue) and classification (orange). For regression, the target/learned functions are scalar output functions, and an image x with digit i is correctly classified if i - 0.5 < f(x) < i + 0.5. It can be seen that both models are in the rich regime, using the minimum number of features.

Finally, the independence from the class labels allows another benefit of being applicable during training. Neural collapse is defined only during the terminal stage of training and assumes 100% classification. A measure without perfect classification can lead to unstable values, typically rising from the pseudo-inverse of  $\Sigma_b$  (See also Appendix H). The NC1 measure shows a dramatic change in logarithmic scale for experiments in Fig. 11(a,b) while our measure is more robust to such differences.

# E CALCULATING THE METRIC AND VISUALIZATION IN PRACTICE

#### E.1 METRIC

In the main text, we used function and operator formalism. In practice, we must use finite approximation. We can use the CKA equation Eq. (6)

$$CKA(A,B) = \frac{\text{Tr}(c(A)c(B))}{\|c(A)\|_F \|c(B)\|_F},$$
(56)

but use  $n \times n$  feature matrices for A and B instead of operators. Let  $A^{1/2} \in \mathbb{R}^{p \times n}$  be feature vectors (of the penultimate layer) of n samples and  $B^{1/2} \in \mathbb{R}^{c \times n}$  be output vectors (of the network) for n samples. We can then approximate CKA by inserting  $A = (A^{1/2})^T A^{1/2}$  and  $B = (B^{1/2})^T B^{1/2}$  into the above equation. Because we already have the square root matrices, the computation cost of CKA is  $\mathcal{O}(npc)$ .

Recall that n can be any number, and can be sampled from both train and test set. As A and B are at most rank p and c respectively, n>p suffices in practice. Because  $p\approx 10^3$  in many models, we can expect  $\mathcal{O}(p^2c)$  computation. In an author's laptop, it can be computed in less than a minute for all experiments. For NTK, p is no longer the width of the last layer, but total number of parameters, quickly becoming infeasible.

#### E.2 VISUALIZATION

To calculate the eigenvalues and eigenfunctions (and thus the quality, utilization, and intensity) for a feature map  $\Phi: \mathcal{X} \to \mathbb{R}^p$ , we need the true input distribution q, which is inaccessible. However, we can use Nyström method (Baker & Taylor, 1979; Williams & Rasmussen, 2006) for a sufficiently large sample size of n>p to approximate the eigenfunctions and eigenvalues of interest. We can create a  $p\times p$  empirical self-covariance matrix  $\Sigma$  from n samples of  $\Phi(x)$ . By diagonalizing  $\Sigma=USU^T$ , we can obtain the empirical eigenvectors and eigenvalues, which can be used to approximate the true eigenvalues and **eigenfunctions**:

$$\rho_k \approx s_k, \qquad e_k(x) \approx \Phi(x)^T u_k / \sqrt{s_k}.$$
(57)

The algorithm is summarized in Algorithm 1.

#### **Algorithm 1** Empirical eigenfunctions and eigenvalues

```
1: \Phi, D_{tr} (Prepare a feature map \Phi and D_{tr}: n samples of training set)
2: \phi \leftarrow \Phi(D_{tr}) (forward transform the input samples to feature vectors)
3: H \in U^T (forward transform singular value decomposition)
```

3:  $U, S, U^T \leftarrow SVD(\phi\phi^T)$  (perform singular value decomposition)

4:  $[u_1, u_2, \dots, u_p] \leftarrow U$  (find the column vectors of U)

5: **for** k in  $1, 2, \ldots, p$  **do** 

6:  $\rho_k \leftarrow s_k$  (approximate eigenvalues)

7:  $e_k(x) \leftarrow \Phi(x)^T u_k / \sqrt{s_k}$  (approximate eigenfunctions)

8: end for

Functions not vectors. Note that eigenfunctions  $e_k's$  are defined beyond the training samples used to compute  $\phi$  and differ from the eigenvectors of the empirical self-covariance matrix. This enables us to evaluate feature quality — via inner products with the target function — on the test set rather than the training set, a capability not available for empirical eigenvectors.

Time complexity. The algorithm is computationally efficient in practice. The main cost in computing Eq. (57) is the SVD, typically with time complexity  $\mathcal{O}(p^3)$ . A good approximation requires only  $n=\mathcal{O}(p)$  datapoints—more than the last layer width. Performing SVD on  $\phi$  instead of  $\phi\phi^T$  yields U and  $\sqrt{S}$  with  $\mathcal{O}(p^3)$  cost. Because last layers are typically narrow (e.g.,  $p\approx 10^3$ ) and the algorithm requires just  $n=\mathcal{O}(p)$  forward passes, the overall computation is lightweight. In an author's laptop, the computation for all experiments takes less than a few minutes.

#### EXAMPLES OF OUR VISUALIZATION METHOD

In this section, we report additional findings made with our visualization tools.

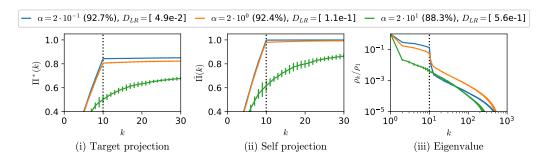


Figure 11: Lazy/rich transition in target downscaling, illustrating how lazier dynamics (ii) use more features and (iii) show a slower decay of eigenvalues (ignoring  $\rho_1$  of constant function, see Appendix D).

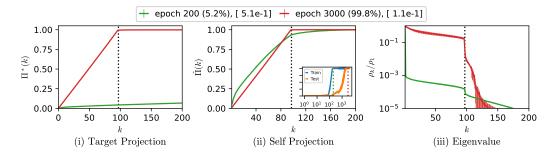


Figure 12: Lazy/rich transition in grokking. A 2-layer transformer is trained on the modular pdivision task. The inset in the middle shows the training and test accuracies, where green and red vertical lines indicate before and after grokking (steps 200 and 3000). Our metric in square brackets shows the transition into the rich regime, while our visualization shows a clear difference in the use of features.

**Low rank bias.** Table 4 shows that low-rank inductive bias is strong in vision tasks, where the number of significant features always equals C (the number of orthogonal functions). See Eq. (8) in Appendix B.4 for the  $D_{eff}(\rho)$  expression. We also observe that our measure  $\mathcal{D}_{LR}$  is small in these setups.

model	dataset	$D_{eff}(\rho)$	$D_{eff}(\hat{Q})$	$\mathcal{D}_{LR}$	test accuracy (%)
ResNet18	MNIST	10.1	10.0	$4.4 \cdot 10^{-3}$	99.7
VGG16	MNIST	10.0	10.0	$1.8 \cdot 10^{-2}$	99.4
ResNet18	CIFAR10	10.4	10.0	$3.1\cdot10^{-3}$	94.8
VGG16	CIFAR10	10.0	10.0	$1.1\cdot 10^{-3}$	93.3
ResNet18	CIFAR100	99.9	99.7	$2.5\cdot10^{-2}$	78.3
VGG16	CIFAR100	95.8	99.4	$7.4\cdot10^{-2}$	71.9

Table 4: Metrics of ResNet18 and VGG16 trained on image datasets

Weight decay is not the source of inductive bias. Fig. 13 shows that weight decay indeed aids low-rank representation and smaller  $\mathcal{D}_{LR}$ . However, the model without any weight decay already shows sufficiently low-rank representations, suggesting that dynamical inductive bias, not the weight decay is the main driving cause of low-rank representations.

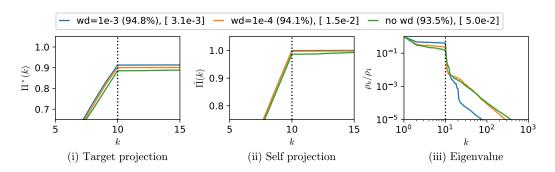


Figure 13: **Effect of weight decay on rich dynamics.** ResNet18 trained on CIFAR10 with varying weight-decay (shown relative value to the fixed learning rate of 0.05). Larger weight decay leads to richer dynamics with smaller  $\mathcal{D}_{LR}$ , but the dynamics is already rich without the weight decay.

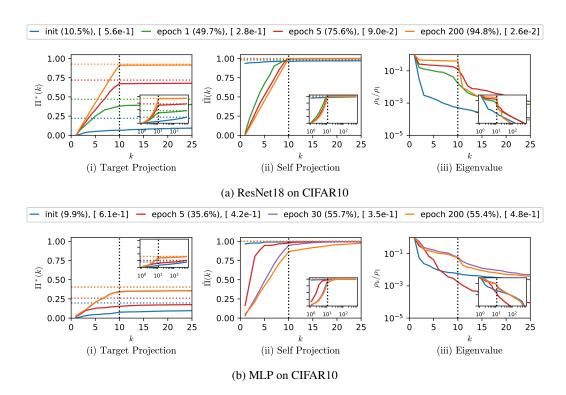


Figure 14: **Role of architecture in the richness of the dynamics.** ResNet18 (a) and a 4-layer MLP of width 512 (b) are trained on CIFAR10, and their metrics are shown at different epochs. ResNet18 concentrates the metrics on the first 10 features after just one epoch, which persists until the end of training. In contrast, MLP shows a less dramatic concentration on the first 10 features and deviates from the rich dynamics around epoch 30 when it begins to use more features.

Role of architecture. As studied in Saxe et al. (2022), the architecture and dataset pair can influence the performance even when dynamics are greedy (race toward shared representation). In Fig. 14, we examine the dynamics of ResNet18 and a 4-layer MLP trained on CIFAR-10. In Fig. 14(a), the training dynamics is restricted to the first 10 significant features, consistent with theoretical work in linear neural networks. In Fig. 14(b), the training is mainly focused on the first 10 features, but we observe that more features are used as training progresses, more similar to the dynamics of linear models. We speculate that the lack of ability to feature learn in the earlier layers leads to the use of additional features and thus leads to lazier training dynamics.

**Rich dynamics**  $\neq$  **generalization.** In Fig. 1, we showed an example where rich dynamics led to poorer performance. Fig. 16 shows a similar experiment in which both rich and lazy models achieve near-perfect performance (> 95% test accuracy). Our visualization clearly shows the different usages of features even when the performances are similar.

**Rich dynamics without underlying data structure.** Fig. 15 illustrates that rich dynamics can occur independently of representation enhancement or underlying data structure. We trained ResNet18 on CIFAR-10 with varying levels of label shuffling and observed that even with fully randomized labels, the model enters the rich regime. This suggests that the dynamical low-rank bias is strong enough to collapse expressive networks into minimal representations, consistent with prior observations of neural collapse under random labeling (Zhu et al., 2021; Hui et al., 2022).

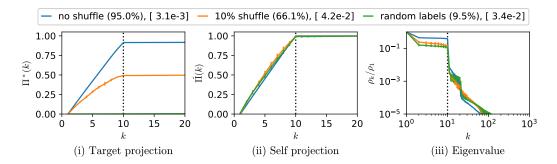


Figure 15: **ResNet18 on 10^4 CIFAR-10 datapoints with shuffled labels.** Models are trained with 0% (blue), 10% (orange), and 100% (green) label shuffling. As shown in the square brackets, all exhibit rich dynamics with  $\mathcal{D}_{LR} < 0.1$ . Visualizations in (ii) and (iii) confirm the use of the top 10 significant features, but varying feature qualities (i) lead to varying test accuracies, shown in parentheses.

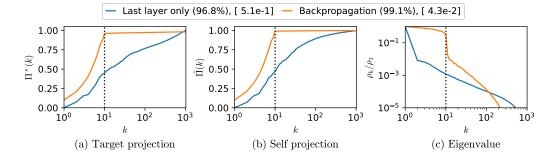


Figure 16: Generalization does not imply rich regime. A CNN with width p=1024 is trained on the full MNIST dataset using two dynamics: last-layer-only training (lazy) and full backpropagation (rich). Although both dynamics achieve test accuracy above 95%, they lead to dramatically different use/significance of features: the metrics concentrate on the first 10 features for the rich regime, and they are more evenly spread for the lazy regime.

**Cross entropy loss.** Fig. 17 shows that our method extends to models trained on cross-entropy loss, while the mathematical interpretation is less straightforward.

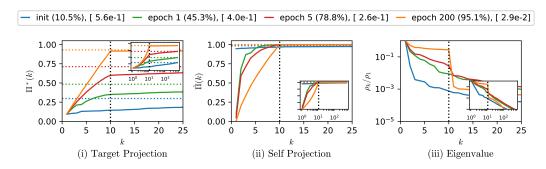


Figure 17: **Training with cross-entropy loss.** We apply our visualization method to a ResNet18 trained on CIFAR-10 with cross-entropy loss. Similar to the MSE-trained model, it achieves small value of  $\mathcal{D}_{LR}$  and shows concentration of feature usage and significance. Insets display all 512 features rather than just the top 25.

G RELATED WORKS ON VISUALIZATION

 Eigendecomposition has long been a standard tool in science and engineering (Golub et al., 1996; Buckley & Xu, 1990; Amadei et al., 1993; te Vrugt & Wittkowski, 2020), demonstrating its value for visualization. Building on this, related empirical works have appeared (Raghu et al., 2017; Morcos et al., 2018; Kornblith et al., 2019). For example, Singular vector canonical correlation analysis (Raghu et al., 2017) resembles eigendecomposition but focuses on performance along leading directions rather than on the number of features used to represent the learned function (i.e., lacking  $\hat{\Pi}$ ). (Morcos et al., 2018; Kornblith et al., 2019) compare different networks, whereas our focus is on before- and after-last-layer activations. Importantly, these approaches were not developed as visualization tools for analyzing rich dynamics.

Similarly, visualizations related to neural collapse use eigendecomposition (Papyan, 2020), but are inherently class-dependent, whereas our metric and visualization are class-independent. See Appendix D.3 and Fig. 10.

The closest line of work is kernel-based visualization (Bordelon et al., 2020; Canatar et al., 2021; Simon et al., 2023), especially (Canatar & Pehlevan, 2022). However, their visualization focuses on cumulative power and the eigenvalue spectrum, whereas ours additionally separates dynamics from performance, providing a complementary perspective. Additionally, the rich dynamics in the literature refers to change in NTK, not the collapse (low-rank bias) of feature kernel.

While visualization is part of our contribution, our main contribution lies in introducing  $\mathcal{D}_{LR}$  as a metric of richness and presenting supporting empirical results. Visualization serves as a complementary tool to illustrate the dynamics and facilitate future studies.

# H EMPIRICAL DETAILS

Here, we share information on our practical setups and their statistical significance. The code can be found at https://anonymous.4open.science/r/RDRFL-0DBA/.

#### H.1 STATISTICAL SIGNIFICANCE

We reproduce tables from the main text but with the addition of one standard deviation values. In the main text, all error bars in the plots represent one standard deviation over at least three runs.

Table 5: Statistical significance of Figure 1(b)

	Full backprop (Rich)	Last layer only (Lazy)
Train acc. ↑	$100 \ (\pm 0)\%$	$98.7\ (\pm0.11)\ \%$
Test acc. ↑	<b>10.0</b> $(\pm 4.3 \cdot 10^{-2})\%$	<b>74.4</b> (±0.19) %
$\mathcal{D}_{LR}\downarrow$	$8.7 \cdot 10^{-3} \ (\pm 4.0 \cdot 10^{-4})$	$6.3 \cdot 10^{-1} \ (\pm 3.2 \cdot 10^{-3})$

Table 6: Statistical significance of Table 1

Epoch	0 (init)	200
Train Acc.↑	$10.0 \ (\pm 0.8)\%$	$10.2(\pm 0.27)\%$
Test Acc.↑	$9.85~(\pm 1.0)\%$	$10.0(\pm 0.33)\%$
${\cal D}_{LR}{\downarrow}$	<b>0.59</b> $(\pm 1.6 \cdot 10^{-2})$	<b>1.0</b> $(\pm 6.4 \cdot 10^{-5})$
$\mathcal{S}_{init} {\downarrow}$	$1.0~(\pm 0)$	$0.20 \ (\pm 1.1 \cdot 10^{-2})$
$\ \theta\ _F^2 \downarrow$	$3.1 \cdot 10^3 \ (\pm 3.5 \cdot 10^2)$	$2.2 \cdot 10^{-5} \ (\pm 2.5 \cdot 10^{-9})$
$NC1\downarrow$	$1.2 \cdot 10^5 \ (\pm 7.2 \cdot 10^2)$	$7.5 \cdot 10^{-14} \ (\pm 1.5 \cdot 10^{-14})$
$Tr(\Sigma_W) \downarrow$	$15.6~(\pm 1.4)$	$7.1 \cdot 10^{-18} \ (\pm 1.1 \cdot 10^{-18})$

Table 7: Statistical significance of Table 2

α	$2 \cdot 10^{-1}$	$2 \cdot 10^0$	$2 \cdot 10^1$
Train Acc.↑	100 (±0)%	100 (±0)%	$100 \ (\pm 0)\%$
Test Acc.↑	92.7 $(\pm 3.4 \cdot 10^{-2})\%$	$92.4 \ (\pm 5.2 \cdot 10^{-2})\%$	$88.3 \ (\pm 3.7 \cdot 10^{-1})\%$
${\cal D}_{LR}\!\!\downarrow$	$4.9 \cdot 10^{-2} \ (\pm 4.4 \cdot 10^{-3})$	$1.1 \cdot 10^{-1} \ (\pm 9.4 \cdot 10^{-3})$	$\mathbf{5.6 \cdot 10^{-1}} \ (\pm 6.1 \cdot 10^{-2})$
$\mathcal{S}_{init}{\downarrow}$	$6.8 \cdot 10^{-2} \ (\pm 2.1 \cdot 10^{-3})$	$4.1 \cdot 10^{-2} \ (\pm 4.5 \cdot 10^{-3})$	$5.2 \cdot 10^{-2} \ (\pm 1.8 \cdot 10^{-2})$
$\ \theta\ _F^2 \downarrow$	$3.4 \cdot 10^3 \ (\pm 1.1 \cdot 10^1)$	$3.2 \cdot 10^3 \ (\pm 4.4)$	$3.2 \cdot 10^3 \ (\pm 2.0)$
$NC1\downarrow$	$2.3 \cdot 10^4 \ (\pm 4.1 \cdot 10^3)$	$3.2 \cdot 10^3 \ (\pm 4.3 \cdot 10^2)$	$8.1 \cdot 10^2 \ (\pm 1.4 \cdot 10^2)$
$Tr(\Sigma_W) \downarrow$	$2.0~(\pm 0.25)$	$3.1 \cdot 10^{-1} \ (\pm 2.3 \cdot 10^{-2})$	$1.2 \cdot 10^{-1} \ (\pm 1.2 \cdot 10^{-2})$

#### H.2 DISCUSSION ON NC1 MEASURE

In Tables 6 and 7, the NC1 measure is numerically unstable due to the pseudo-inverse of  $\Sigma_b$ . For stability, We used  $\Sigma_b^{\dagger} \approx (\Sigma_b + 10^{-4} I)^{\dagger}$ . In addition to NC1, we consider a similar measure  $Tr(\Sigma_W)$  Papyan et al. (2020) for collapse of within class variability, which is also reported in Tables 6 and 7. The measure  $Tr(\Sigma_W)$  also shows a similar trend to NC1.

#### H.3 DATASET DETAILS

We use publicly available datasets, including MNIST Deng (2012) and CIFAR-10/100 Krizhevsky et al. (2009) from Pytorch Paszke et al. (2019), and mod p division task from https://github.com/teddykoker/grokking.

Table 8: Statistical significance of Table 3

Task	Architecture	Condition	$D_{LR}$
Mod	2-layer	Step 200 (before grokking)	$0.51\ (\pm0.04)$
97	transformer	Step 3000 (after grokking)	$0.11\ (\pm0.01)$
CIEAD		learning rate = 0.005	$0.053\ (\pm0.01)$
CIFAR- 100	ResNet18	learning rate = $0.05$	$0.025\ (\pm0.004)$
100		learning rate = $0.2$	$0.039\ (\pm0.007)$
CIFAR-	ResNet18	weight decay = 0	$0.05\ (\pm0.003)$
10	Resiletto	weight decay = $10^{-4}$	$0.015\ (\pm0.003)$
		weight decay = $10^{-3}$	$0.003\ (\pm0.001)$
CIFAR-	ResNet18		$0.026\ (\pm0.003)$
10	MLP	•	$0.48\ (\pm0.07)$
CIFAR-		no label shuffling	0.031 (±0.004)
10	ResNet18	10% label shuffling	$0.042\ (\pm0.003)$
10		full label shuffling	$0.034\ (\pm0.004)$
MNIST	CNN	full backpropagation	$0.043\ (\pm0.005)$
MINIST	CIVIV	last layer only training	$0.51\ (\pm0.008)$
CIFAR-	VGG-16	without batch nomralization	$0.66 (\pm 0.01)$
100	¥ GG-10	with batch normalization	$0.073\ (\pm0.002)$

#### H.4 MODEL DETAILS

Our model implementations are based on publicly available code assets, including VGG16 from PyTorch Paszke et al. (2019), ResNet18 from Nakkiran et al. Nakkiran et al. (2021), and the modular division task from https://github.com/teddykoker/grokking.

As described in the main text, we only trained ours models on MSE loss where the targets are one-hot vectors (up to scaling constant). The constant  $\alpha$  is used to scale the output  $y \to y/\alpha$ .

Table 9: Dataset details

Figure	dataset	training sample count	α	batch size
Fig. 1	Encoded MNIST (Appendix H.4.1)	60,000 (all)	1/3	128
Table 2	MNIST	60,000 (all)	1	128
Fig. 2	CIFAR-10	•	1/3	128
Fig. 3	CIFAR-100	50,000 (all)	1/10	128
Fig. 5	CIFAR-100	50,000 (all)	1/10	128
Fig. 11	MNIST	1,000		128
Fig. 12	mod-p division	4,656	1	512
Fig. 15	CIFAR-10	10,000	1/3	128

# H.4.1 ENCODED MNIST

In Fig. 1, we encoded the labels as one-hot vectors on the first 10 pixels of the MNIST dataset. For the training set, we encoded the true labels, but for the test set, we encoded random labels.

#### H.4.2 CIFAR-10/100

For CIFAR-10/100, we use standard augmentation using randomcrop 32 with padding 4 and random horizontal flip with probability 0.5. We also use standard normalization with mean [0.4914, 0.4822, 0.4465] and standard deviation [0.2023, 0.1994, 0.2010] for each channels.

Table 10: Training details

Figure	optimizer	learning rate	momentum / beta	weight decay	epochs	learning rate scheduling
Fig. 1	Adam	$1\cdot 10^{-3}$	[0.9, 0.999]	0	100	None
Table 1	SGD	$1 \cdot 10^{-6}$	0	$1 \cdot 10^{-3}$	200	None
Fig. 2	SGD	$5\cdot 10^{-2}$	0.9	$5\cdot 10^{-5}$	200	$\times 0.2$ per 60 epochs
Fig. 3	SGD	$1 \cdot 10^{-2}$	0.9	$1\cdot 10^{-5}$	400	None
Fig. 5	SGD	$5\cdot 10^{-2}$	0.9	$5\cdot 10^{-5}$	400	$\times 0.2$ per 60 epochs
Fig. 11	Adam	$1\cdot 10^{-3}$	[0.9, 0.999]	0	100	None
Fig. 12	Adam	$1\cdot 10^{-3}$	[0.9, 0.98]	0	4000	None
Fig. 15	SGD	$5\cdot 10^{-2}$	0.9	$5\cdot 10^{-5}$	200	$\times 0.2$ per 60 epochs

#### H.5 COMPUTE RESOURCES

The models ran on GPU cluster containing RTX 1080 (8GB), RTX 2080 (8GB), RTX3060 (12 GB), and RTX3090 (24GB). The typical time to train a model is 2 hours, but it varies from 10 minutes to 6 hours depending on the experiment. The evaluation metrics take less than 5 minutes and may need up to 2GB of CPU memory.

# H.6 USE OF LLMS

Large language models (LLMs) were used to polish writing to make the paragraphs more concise.