Approaching Deep Learning through the Spectral Dynamics of Weights

Anonymous authors

Paper under double-blind review

ABSTRACT

We propose an empirical approach centered on the spectral dynamics of weights the behavior of singular values and vectors during optimization—to unify and clarify several phenomena in deep learning. We identify a consistent bias in optimization across various experiments, from small-scale "grokking" to large-scale tasks like image classification with ConvNets, image generation with UNets, speech recognition with LSTMs, and language modeling with Transformers. We also demonstrate that weight decay enhances this bias beyond its role as a norm regularizer, even in practical systems. Moreover, we show that these spectral dynamics distinguish memorizing networks from generalizing ones, offering a novel perspective on this longstanding conundrum. Additionally, we leverage spectral dynamics to explore the emergence of well-performing sparse subnetworks (lottery tickets) and the structure of the loss surface through linear mode connectivity. Our findings suggest that spectral dynamics provide a coherent framework to better understand the behavior of neural networks across diverse settings.

024 025 026

027

004

006

008 009

010 011

012

013

014

015

016

017

018

019

021

1 INTRODUCTION

Interest in neural networks has exploded in the past decade. Capabilities are rapidly improving, and deployment is ever-increasing. Yet, although issues with these technologies now have social repercussions (Bender et al., 2021; Bommasani et al., 2021), many fundamental questions regarding their behavior remain unanswered.

For instance, despite extensive research, we still lack a complete understanding of the implicit biases of neural networks trained via stochastic optimization (Neyshabur et al., 2014). Even basic questions regarding the role of regularization like weight decay (Hanson & Pratt, 1988; Krogh & Hertz, 1991; Zhang et al., 2018a) have only partial answers (Van Laarhoven, 2017; Andriushchenko et al., 2023; Yaras et al., 2023b). Perhaps most vexing, we lack a complete explanation for how neural networks generalize, despite having the capacity to perfectly memorize the training data (Zhang et al., 2021). Such an explanation may allow us to design better algorithms, however a lack of understanding makes the deployment of neural networks vulnerable to uninterpretable errors across fields (Szegedy et al., 2013; Ilyas et al., 2019; Hendrycks et al., 2021; Zou et al., 2023).

041 Although theoretical explanations have been put forward, these studies are often limited to special 042 settings like deep linear networks (Arora et al., 2018; 2019) or infinite-width systems (Jacot et al., 043 2018), and arguments may rely on unsubstantiated or impractical assumptions like near-zero initial-044 ization. On the empirical side, a growing body of work in interpretability has attempted to reverseengineer neural networks (Rahaman et al., 2019; Barak et al., 2022; Nanda et al., 2023), but given the difficulty of the task, the systems of interest have been very small-scale, and the methodology 046 for analysis quite bespoke and challenging to scale. A third category of work aims at understanding 047 empirical behavior from a higher level (Zhang et al., 2021; Huh et al., 2022; Yu & Wu, 2023), but 048 while these works often study larger-scale systems, they often focus on more abstract objects like 049 the gram matrix (Huh et al., 2022) or the Neural tangent kernel (NTK) (Fort et al., 2020), and thus 050 do not have the granularity and predictive power of the previous two categories. 051

To bridge these gaps, we propose a task-agnostic, unifying perspective of many disparate phenomena
 in deep learning across many different practical tasks and architectures, including image classifica tion with ConvNets, image generation with UNets, speech recognition with LSTMs and language

068

069

071 072



Figure 1: (a) Schematic for the spectral dynamics of a weight matrix. As training proceeds, top 064 singular vectors become stable and top singular values grow disproportionately large. (b) Singular value evolution for a single matrix in a Transformer, where each line is a single singular value and 065 color represents rank. We see a disproportionate trend where large singular values grow larger faster. 066 (c) Previous works have used SV basis alignment between layers to prove similar theoretical results, however actual alignment between consecutive layers is not strong, which we describe in Section 4. We explore these spectral dynamics of weights and connect them to generalization, regularization, and seemingly unrelated phenomena like linear mode connectivity. 070

073 modeling with Transformers. Through extensive experiments, we examine the dynamics (i.e., evolu-074 tion over training) of singular values and singular vectors of weight matrices; the spectral dynamics 075 of weights. We observe a few key properties: singular values evolve unequally, with larger ones evolving faster, as a result top singular vectors stabilize toward the end of training, and for top ranks 076 we see alignment between neighboring layers' singular vectors, though this varies somewhat with 077 the setting. We preview these in Figure 1. We are motivated to study these dynamics specifically as optimization is posited to be one of the fundamental process driving deep learning (Nagarajan & 079 Kolter, 2019; Zhang et al., 2021), and the SVD is fundamental to every matrix. We detail how these 080 properties connect to generalization and many other phenomena in the following paragraphs. 081

Contributions: As a test bed for understanding generalization, Power et al. (2022) introduce the 083 "grokking" phenomenon, where a small-scale model on arithmetic tasks initially minimizes the training loss but performs poorly on validation data, then with much more training suddenly mini-084 mizes the validation loss. In particular, Nanda et al. (2023) showed that in modular arithmetic, the 085 feature learning that occurs during grokking could be completely reverse-engineered from the final 086 weight matrices. Although this description is precise, it is limited to these particular tasks. In Sec-087 tion 3, we notice a task-agnostic view of grokking, observing that the drop in validation loss during 880 grokking coincides with the simultaneous discovery of low-rank solutions across all weight matrices 089 in the network, whether it be modular arithmetic or an image-classification setting. We also find that this transition relies on weight decay, echoing existing works (Lyu et al., 2023; Liu et al., 2023) as 091 neither grokking nor low-rank weights occur without sufficient weight decay.

092 Though the common tie between low-rank weights and generalization in grokking, suggests a con-093 nection between rank and generalization, grokking is typically studied on synthetic tasks with very 094 small-scale models like single-layer Transformers or small MLPs, and requires very particular hy-095 perparameter settings (Power et al., 2022; Nanda et al., 2023; Gromov, 2023; Kumar et al., 2023). If 096 our perspective is to be useful, it needs to scale to larger systems. Thus, we turn to common empirical tasks drawn from the literature like image classification, image generation, speech recognition 098 and language modeling as well as varied and larger networks like VGG (Simonyan & Zisserman, 2014), UNet (Ronneberger et al., 2015), LSTM (Hochreiter & Schmidhuber, 1997b) and multi-layer 099 Transformers (Vaswani et al., 2017). 100

101 In Section 4, we demonstrate that the spectral dynamics are biased toward effective rank minimiza-102 tion across various practical neural networks in complex settings. Although this behavior echoes 103 theoretical predictions in the deep linear setting, we find that the behavior of networks disagrees 104 with a common theoretical assumption about low-rank dynamics: the alignment of singular vectors 105 in consecutive layers (Saxe et al., 2014; Arora et al., 2018; 2019; Milanesi et al., 2021). Thus, the rank minimization mechanism may differ from what the theory describes. It is notable too that our 106 hyperparameter settings are drawn from existing literature, thus the trend toward rank minimization 107 coincides with well-generalizing networks across settings.

108 Given one particularly notable ingredient for grokking was a very high level of weight decay, in 109 Section 5 we empirically connect rank minimization to weight decay, showing that weight decay 110 promotes rank minimization across architectures and tasks, echoing Section 3. In addition, in some 111 cases, it also appears to promote singular vector alignment in consecutive weights despite the nonlin-112 earities between layers. Although weight decay explicitly penalizes norm, studying spectral dynamics allows us to observe a host of effects including on rank and alignment. Such effects may help 113 in understanding the reason weight decay is useful, as norm-based explanations have been found 114 insufficient (Andriushchenko et al., 2023). 115

116 To test the explanatory power of the framework, we turn to the classic memorization experiments of 117 Zhang et al. (2021), who demonstrated that even small networks can memorize random labels, thus 118 any arguments about generalization cannot be capacity-based alone. In Section 6 and Appendix B, we show that training with random labels leads to high-rank solutions, while rank with true labels is 119 much lower. We also find that while random labels do not align consecutive layers, true labels do, 120 which is surprising given the non-linearities between layers. This echoes prior discussion on rank 121 and generalization. Through spectral dynamics, we see a clear distinction between generalizing and 122 memorizing networks, which provides a foothold toward better theoretical understanding. 123

Our results suggest that viewing neural networks through the lens of spectral dynamics can shed light 124 125 on several generalization-related phenomena, but we suspect there are broader connections. In the literature, many curious and unexplained phenomena regarding neural networks exist. We take two 126 as case studies. First, the lottery ticket hypothesis (LTH) (Frankle & Carbin, 2018) and second, linear 127 mode connectivity (LMC) (Nagarajan & Kolter, 2019; Frankle et al., 2020; Neyshabur et al., 2020). 128 We find that global magnitude pruning, a standard procedure for finding lottery tickets, preserves top 129 singular vectors and acts like a low-rank pruning. We also see that the ability to interpolate between 130 models in LMC strongly correlates with sharing top singular vectors. With these results, we note 131 that the two phenomena can be seen as aspects of the spectral dynamics of weights, bringing them 132 under the umbrella of prior sections. For detailed discussion see Section 6 and Appendix C.

- To summarize the discussion above, by studying the spectral dynamics of weights, we find:
 - Grokking is intimately linked to rank minimization;
 - Rank minimization is a general phenomenon in more complex tasks;
 - Weight decay acts implicitly as a low-rank regularizer;
 - Generalizing solutions have a lower rank than memorizing ones; and
 - Top singular vectors are preserved when performing magnitude pruning and while linearly interpolating between connected modes.

All of these phenomena and effects have previously been studied in isolation to varying degrees, but by approaching deep learning through spectral dynamics, we aim to provide a common language for probing and understanding neural networks. Code for all experiments will be released.

145 146 147

133

135

136

137 138

139

140

141

142 143

144

2 RELATED WORK

148 149

150 2.1 SINGULAR VALUE DYNAMICS

151 Prior work on deep linear networks (Arora et al., 2019; Milanesi et al., 2021) suggests that rank 152 minimization may better describe implicit regularization in deep matrix factorization than simple 153 matrix norms. See Arora et al. (2018) (Appendix A) for a detailed argument. However, a critical 154 assumption in these works is "balanced initialization." This means that for consecutive matrices W_i and W_{i+1} in the product matrix $\prod_j W_j$, we have $W_{i+1}^\top W_{i+1} = W_i W_i^\top$ at initialization. Decompos-155 156 ing these matrices with SVDs and leveraging orthogonality leads to matching left and right singular 157 vectors between consecutive matrices. See Appendix A for a detailed explanation. Consequently, 158 the product of the diagonals will evolve in a closed-form manner, with larger singular values grow-159 ing faster than smaller ones. As shown by Arora et al. (2019), this translates to rank-minimizing behavior with increasing depth in the matrix products. This formula is also empirically validated 160 for linear matrix factorization problems. Similar results have been derived for tensor products and 161 other structured settings (Saxe et al., 2014; Yaras et al., 2023a). (Ji & Telgarsky, 2019) show that alignment between layers will happen specifically for deep linear networks with infinite training.
 Still, there is no reason to believe standard networks obey this balancedness condition under practical initialization procedures. In Section 4, we explore how these conclusions and assumptions hold for much larger, practical neural networks that are far from linear.

166 167

168

2.2 LOW-RANK PROPERTIES

Another line of research focuses on more general low-rank biases. Early work explored norms as an implicit bias (Gunasekar et al., 2017). Theoretical analyses reveal that norms or closed-form func-170 tions of weights might be insufficient to explain implicit regularization, but they do not necessarily 171 contradict the possibility of rank minimization (Razin & Cohen, 2020; Vardi & Shamir, 2021). Nu-172 merous studies investigate low-rank biases in various matrices, including the Jacobian (Pennington 173 et al., 2018), weight matrices (Le & Jegelka, 2021; Martin & Mahoney, 2020; 2021; Frei et al., 174 2022; Ongie & Willett, 2022), Gram matrix (Huh et al., 2022), and features (Yu & Wu, 2023; Feng 175 et al., 2022). Additionally, research suggests that dynamics influence the decay of rank (Li et al., 176 2020; Chen et al., 2023; Wang & Jacot, 2023). Orthogonally, weight decay has a long history 177 as a regularizer explicitly penalizing parameter norm, which can be used for norm-based general-178 ization bounds (Bartlett, 1996), but these bounds do not seem to explain the success of practical 179 systems (Nagarajan & Kolter, 2019; Jiang et al., 2019). Some works establish connections between weight decay and rank minimization in idealized settings (Ziyin et al., 2022; Galanti et al., 2022; 180 Zangrando et al., 2024; Ergen & Pilanci, 2023; Parhi & Nowak, 2023; Shenouda et al., 2023), which 181 may be connected to generalization (Razin & Cohen, 2020). We are particularly interested in how 182 far these connections extend in practice. 183

184 185

3 GROKKING AND RANK MINIMIZATION

Power et al. (2022) first noticed a surprising phenomenon they called "grokking" where models quickly fit the training data on toy tasks, then after a long period of training, very quickly generalize on the validation data. Later, others found that this phenomenon can occur in a relaxed fashion (Thi-lak et al., 2022) on very simple models and different datasets (Liu et al., 2022; Gromov, 2023; Kumar et al., 2023; Xu et al., 2023), and that weight decay seems critical to cause it (Lyu et al., 2023; Liu et al., 2023; Tan & Huang, 2023).

Motivated by theoretical work that proposes connections between rank and generalization (Razin & Cohen, 2020), weight decay and rank (Galanti et al., 2022; Timor et al., 2023; Yaras et al., 2023b;
Zangrando et al., 2024), and the importance of weight decay for grokking (Power et al., 2022; Lyu et al., 2023; Liu et al., 2023) in simple settings, we evaluate the potential connection between rank and grokking in neural networks. This offers a complementary perspective on grokking with other descriptions such as Fourier decomposition (Nanda et al., 2023), the simplification of linear decision boundaries (Humayun et al., 2024), the connection to double descent (Davies et al., 2022), and the discovery of a sparse solution (Merrill et al., 2023).

We replicate grokking in two settings: a single-layer Transformer for modular addition (Nanda et al., 2023), and a 12-layer MLP for MNIST image classification (Fan et al., 2024) (see Appendix D for details). Inspired by work in the deep linear case (Saxe et al., 2014; Arora et al., 2019; Milanesi et al., 2021; Yaras et al., 2023b), we track the evolution of singular values for individual weight matrices. To gain a high-level overview of all parameter evolutions, we compute the (normalized) effective rank of a matrix W (Roy & Vetterli, 2007) with rank R as

207 208

$$\operatorname{EffRank}(W) := -\sum_{i=1}^{R} \frac{\sigma_i}{\sum_j \sigma_j} \log \frac{\sigma_i}{\sum_j \sigma_j} , \qquad (1)$$

209 210 $NormEffRank(W) := \frac{EffRank(W)}{R} , \qquad (2)$

where σ_i 's are the singular values of matrix W and EffRank(W) is the entropy of the normalized singular value distribution. As the probability mass concentrates, the effective rank decreases. We plot NormEffRank(W) to compare across layers and time.

In addition, inspired by the assumptions of balancedness made by prior work (Arora et al., 2018; 2019), we examine the alignment of consecutive weight matrices in the Transformer. To examine



Figure 2: Grokking and Spectral Dynamics. Top row: Grokking for Transformers on modular 235 addition (Nanda et al., 2023). Bottom row: Grokking for a 12-layer MLP on MNIST (Fan et al., 236 2024). 1st column: Training and validation error. 2nd column: A visualization of singular value 237 evolution for the first attention parameter and the second MLP layer, where each line represents 238 a single singular value and the color represents the rank. **3rd column:** Effective rank of all layers 239 (Eqn. 1). 4th column: A visualization of the alignment (Eqn. 3) between the embedding and the first attention parameter, and the first and second MLP layers, where the y-axis corresponds to index i of 240 the diagonal. We see that grokking co-occurs with a transition to low-rank weights. In addition, there 241 is an alignment that begins early in training that evolves up the diagonal. In the image classification 242 case, we see a similar rank transition, though alignment appears seemingly out of nowhere. 243

and quantify this alignment between SVDs of consecutive matrices in a network at training time t, i.e.,

$$W_{i} = \sum_{j=1}^{R} \sigma_{j}(t) u_{j}(t) v_{j}(t)^{\top}, \qquad W_{i+1} = \sum_{k=1}^{R} \sigma_{k}'(t) u_{k}'(t) v_{k}'(t)^{\top} ,$$

we compute,

244 245

246

252 253 254

255

256

257

$$A(t)_{jk} = |\langle u_j(t), v'_k(t) \rangle| \quad , \tag{3}$$

where the absolute value is taken to ignore sign flips in the SVD computation. We then plot the diagonal of this matrix $A(t)_{ii} \forall i \leq 100$ over time. For exact details on how alignment is computed for different architectures and layers more complex than the fully connected case, see Appendix D.

In Figure 2, we see a tight connection: the sudden drop in validation loss coincides precisely with 258 the onset of low-rank behavior in the singular values. Examining inter-layer alignment during train-259 ing, we observe that the final low-rank solution gradually emerges from the model's middle ranks. 260 Conversely, in Figure 19, the grokking phenomenon is absent without weight decay, and no low-261 rank solution seems to develop. Additionally, when using 90% of the data and no weight decay, 262 generalization still coincides with effective rank minimization. Fan et al. (2024) noted that in deep 263 MLPs, grokking coincided with a feature rank decrease, which we show stems from the parameter 264 rank decrease here. The familiar reader will also note that Nanda et al. (2023) previously showed 265 that the particular solution found in modular addition is a low rank fourier decomposition, so our 266 observations on low rank weights will follow, yet the same structure also applies to the MLP where 267 such reverse-engineering is difficult. In the following sections we argue that rank minimization is a perspective that can apply in more complex settings when one does not know what to look for in the 268 weights, and it may be possible to interpret the neural network via the top ranks (Praggastis et al., 269 2022).

270 4 SPECTRAL DYNAMICS ACROSS TASKS 271

Inspired by the results on grokking and prior work on deep linear networks that studies the evolution of the SVD of the weight matrices (Saxe et al., 2014; Arora et al., 2018; 2019; Milanesi et al., 2021; Yaras et al., 2023a), we apply the same analysis to larger, more practical systems. We show that the trends we saw in the analysis of grokking mostly hold true across networks and tasks at a much larger scale, even though our findings do occasionally deviate from theoretical predictions.

4.1 Methodology

Cun, 1998);

Our experiments aim to examine reasonably sized neural networks across a variety of tasks. We select models and tasks that are representative of current applications. Specifically, we focus on:

• Image classification with CNNs (VGG-16 (Simonyan & Zisserman, 2014)) on CI-FAR10 (Krizhevsky, 2009);

• Image generation through diffusion with UNets (Ronneberger et al., 2015) on MNIST (Le-

284 285

277 278

279 280

281 282

283

- 286
- 287 288

289

290

291

- Speech recognition with LSTMs (Hochreiter & Schmidhuber, 1997b) on LibriSpeech (Panayotov et al., 2015); and
- Language modeling with Transformers (Vaswani et al., 2017) on Wikitext-103 (Merity et al., 2016).

Training hundreds of runs for each of the above experiments is computationally expensive, limiting the scale of models we can explore. We primarily adopt hyperparameters from existing literature, with minor modifications for simplicity. This ensures that any correlations observed are likely a reflection of common practices, not introduced bias on our part. We also provide evidence with larger scale (up to 3B parameters) in Appendix D.5, from the Pythia suite (Biderman et al., 2023).

297 The primary evidence in this section comes from computing the SVDs of weight matrices within 298 the models. Consequently, we disregard 1D bias and normalization parameters in our analysis. 299 Indeed previous research suggests that in some cases these parameters are not crucial for performance (Zhang et al., 2018b; Mohan et al., 2019; Karras et al., 2023). Due to the large number of 300 matrices in these models, we present plots of individual layers' matrix parameters and statistics that 301 summarize behavior across layers for conciseness of presentation. Hundreds of thousands of plots 302 were generated for this study, making it impossible to include them all. Full experimental details, 303 including the choice of hyperparameters, are available in Appendix D. 304

305 306

4.2 EFFECTIVE RANK MINIMIZATION

Building on theoretical (Saxe et al., 2014; Arora et al., 2019; Milanesi et al., 2021; Boix-Adserà et al., 2023; Yaras et al., 2023a) and empirical (Dittmer et al., 2019; Martin & Mahoney, 2020; 2021; Boix-Adserà et al., 2023) findings, we investigate effective rank minimization across parameters in larger models and on a more diverse variety of tasks. Figure 3 reveals a consistent trend: the effective rank of network parameters generally decreases throughout training, regardless of the specific parameter or network architecture. This suggests a progressive "simplification" of the network as training progresses.

314 We further conduct a singular-value pruning experiment to explore the relationship between low-315 rank behavior and model performance. We prune either the top or bottom half of the singular values 316 for each weight matrix in the network and then evaluate the pruned model at each training step. 317 Given their importance in \mathcal{L}^2 space, we expect the top singular values to capture the most critical 318 information for the network's function. Figure 4 confirms this, demonstrating that the pruned pa-319 rameters, without further training, can closely approximate the full model's performance. It is not 320 necessarily obvious that pruning would have this effect. In particular, simultaneously pruning lower 321 components across all layers may lead to losing some critical signal that must be passed between layers, or it could be that small-magnitude singular values may provide some important regularizing 322 noise. In later sections, we will rely on this observation that large singular values are more critical 323 to the function of the network.



Figure 3: Top row: Singular value evolution for a single matrix in the middle of each model. Each line represents a singular value, while color represents rank. Notice the unequal evolution where top singular values grow at a disproportionate rate. Bottom row: Normalized effective rank (Eqn. 1) evolution visualized in color for different matrices across architectures and time. As we move down the y-axis, the depth of the parameters in the model increases, while the x-axis tracks training time. Notice decreasing effective rank across nearly all parameters, though the magnitude differs across layers. The block-like patterns for VGG are likely due to different channel dimension sizes. The banding in the UNet, LSTM, and Transformer is due to the differences between convolutional and linear layers, residual block connections, and attention and fully connected layers, respectively. The sharp transition midway through training in the VGG case is likely due to a $10 \times$ learning rate decay.



Figure 4: **Left plot:** Training loss. **Right plot:** Validation loss. **Red** is the full model. **Blue** is post-training pruning the bottom half of the SVD for every matrix in the model that is not the final layer. Green is post-training pruning the top half of the SVD. Notice that for all models, keeping the top half of the SVD is close to the full model performance, supporting the idea that the top directions provide a better approximation to the function.

4.3 ALIGNMENT OF SINGULAR VECTORS BETWEEN LAYERS

Similar to the analysis of grokking, we investigate the alignment between consecutive layers in the larger neural networks considered in this section. We not only employ the alignment matrix defined in Eqn. 3 but also derive and plot a scalar measure for alignment based on the top diagonal entries:

$$a(t) = \frac{1}{10} \sum_{i=1}^{10} A(t)_{ii}$$
(4)

For specific details on calculating this measure in diverse architectures and complex layers (beyond fully connected layers), please refer to Appendix D.

Figure 5 reveals a key finding: the theoretical assumption of **balanced initialization**, which posits aligned singular value decompositions (SVDs) between weight matrices (Arora et al., 2018; Saxe et al., 2014), does not hold true at the start of training in these larger networks. Additionally, unlike the linear case discussed in Du et al. (2018), the alignment does not appear to remain static through-



Figure 5: Neighboring layer alignment of singular vectors. Left plot: The diagonal of the alignment matrix $A(t)_{ii}$ (Eqn. 3) vs. training time for a single pair of matrices in the middle of each model. We see a small amount of alignment in the top ranks between layers shortly after training begins, but this becomes more diffuse over time. **Right plot:** Alignment metric (Eqn. 4) for pairs of matrices for depth vs. training time. It is hard to make out a global trend across models, though the LSTM shows a weak signal around Epoch 1 when the initial alignment occurs, and the Transformer case has a banding pattern with depth due to alignment between the query and key matrices that have no nonlinearity in between.

out training. However, a weak signal of alignment in the top ranks develops and disappears. This trend is somewhat reminiscent of the theoretical result provided by Mulayoff & Michaeli (2020) for linear networks under the assumption of whitened input data. Still, the weakness of the observed signal means that existing theoretical models do not capture the complexities of neural network training.

5 THE EFFECT OF WEIGHT DECAY

392

394

395

396

397

398 399

400

401 In light of the previously observed evolution of singular values, we investigate a proposed effect 402 of weight decay. Though weight decay explicitly penalizes the norm of weights, there is evidence 403 that complicates the connection between the norm and generalization for neural networks (Razin 404 & Cohen, 2020; Andriushchenko et al., 2023), meaning we do not have a full understanding as to 405 why weight decay may be useful. Alternatively, some theoretical (Boix-Adserà et al., 2023; Razin 406 & Cohen, 2020; Yaras et al., 2023a; Timor et al., 2023; Ongie & Willett, 2022; Galanti et al., 2022; 407 Zangrando et al., 2024) and empirical works (Galanti et al., 2022; Boix-Adserà et al., 2023) propose 408 a connection with the rank of matrices in constrained settings. Still, a comprehensive connection to 409 larger empirical networks has not yet been demonstrated.

410 We speculate on the intuition of the mechanism in more practical settings. In its simplest form, 411 weight decay involves the optimization arg min_W $\mathcal{L}(W) + \lambda \|W\|_F^2$, where $\|W\|_F^2 = \sum_{i=1}^R \sigma_i^2$ 412 with singular values σ_i of weight matrix \tilde{W} with rank \hat{R} . We saw previously that larger singular 413 values of neural networks grow faster (Fig. 3, top row) and that the top singular vectors are much 414 more useful for minimizing task loss than the bottom ones (Fig. 4). Thus, with minor weight decay 415 regularization, one straightforward solution for the network may be to minimize the rank of a given 416 weight matrix while preserving the top singular values to minimize $\mathcal{L}(W)$. Timor et al. (2023) argue 417 a similar effect: if all singular values are less than one, the norm of activations will shrink with depth, so it will be impossible to pass signals from input to output in sufficiently deep networks with even 418 penalization. Thus it is better for a few singular values to be sufficiently large, while the rest can be 419 very small. 420

Figure 6 shows that adding weight decay produces this exact low-rank behavior, while too much
weight decay leads to complete norm collapse. The exact choice of "too much" varies across architectures and tasks.

424 Despite the low-rank regularization, we do not see particularly tight alignment in the top singular 425 vectors, with the exception of the highest weight decay Trasnformer (Figure 17). This behavior is 426 quite reminiscent of the balancedness condition (Arora et al., 2018; 2019; Du et al., 2018), though the 427 Transformer considered here has nonlinearities and much more complex structure. It is curious that 428 the trend reverses for only this architecture. We also provide additional evidence in Appendix D, where Figure 18 shows that the solutions with very high weight decay are still performant, even 429 though they are much lower rank. Though it is difficult to argue as simple a trend as "lower rank 430 equals better generalization" because one does not know the minimal rank necessary for a given task, 431 we note that the role of weight decay for improving generalization is tied up with its function as a



Figure 6: SV evolution for a single matrix and normalized effective rank (Eqn. 1) across matrices over time, where the rows use differing amounts of weight decay. From top to bottom, for VGG we use coefficients {0,0.01,0.01,0.1}, while for other networks we use coefficients {0,0.1,1,10}. Higher weight decay coefficients promote more aggressive rank minimization. VGG uses SGD w/ momentum, while the rest use AdamW (Loshchilov & Hutter, 2017), which may explain the earlier norm collapse.



Figure 7: Top row: results with true labels. Bottom row: results with random labels. We see that
the middle layers have a lower effective rank when using true labels and that alignment in the middle
layers persists throughout training, unlike in the random label case. We emphasize this alignment
occurs despite the nonlinearities.

rank regularizer. In addition, although we lack precise tools to entirely interpret complex models, when there are only a few ranks per matrix, it may become possible to extend analysis efforts (Nanda et al., 2023; Praggastis et al., 2022) to more complex domains.

6 ADDITIONAL CONNECTIONS

Here we briefly preview some connections between spectral dynamics and additional phenomena.

Memorization vs. Generalization: In Figure 7, we replicate the core memorization experiment of Zhang et al. (2021), which highlighted the ability of modern neural networks to memorize perfectly

486 even random labels across several supervised learning tasks. We find that when training with random 487 labels, networks with higher rank final parameters are obtained as opposed to when training with 488 true labels. Thus, the spectral dynamics can distinguish between memorization (of random labels) 489 and generalization. We also see an alignment structure between the middle layers that disappears 490 with random labels, perhaps as it is necessary in order for the network to pass signals from input to output. This echoes the pattern of grokking in Section 3, where generalization came with a 491 transition to low-rank and alignment in middle layers. We expand this experiment to more settings 492 in Appendix **B** with additional figures. 493

494 Lottery Tickets: On very small networks, Frankle & Carbin (2018) found the existence of sparse 495 sub-networks via magnitude pruning, keeping only the top p% of weights globally by magnitude, 496 that could train to similar performance as the full network. For larger image classification networks, Frankle et al. (2020) observed that in order to find such sparse subnetworks, it was necessary to train 497 till the end in order to acquire the pruning mask that could be used retroactively in training. This cu-498 rious observation still lacks a compelling explanation. We show that such global magnitude pruning 499 functions similarly to low-rank pruning, thus the lottery ticket masks found by rewinding (Frankle 500 et al., 2020), are effectively low-rank masks for the singular components that will become important 501 at the end of training. Training the masked network leads to the same dynamics in these compo-502 nents. However, taking masks from too early in training leads to poor approximation of these final 503 components and simultaneously stunts training. We provide detail on this discussion in Appendix C. 504

Linear Mode Connectivity (LMC): Linear Mode Connectivity (Nagarajan & Kolter, 2019; Frankle 505 et al., 2020) refers to the property that models that share a portion of the training trajectory can be 506 averaged in weight-space to yield a stronger model (Wortsman et al., 2022; Ramesh et al., 2022). 507 This phenomenon indicates that, after some training, the loss surface is quite convex in a subspace, 508 even though the optimization problem is theoretically highly non-convex. As all fine-tuning from 509 pre-trained models stays in this convex space (Neyshabur et al., 2020; Li et al., 2022; Sadrtdinov 510 et al., 2023), an explanation for what underlies LMC could help to clarify the role of pre-training, and 511 may lead to faster fine-tuning. We show that LMC is tied with singular vector sharing. In particular, 512 when models can be averaged they share top singular vectors between weights, and when they cannot 513 they do not. This is an outcome of the early stability of top singular vectors, which arises due to the unequal evolution of singular values. It is also straightforward to explain the large euclidean 514 distance between checkpoints that can be averaged, as they only share a very small portion of the 515 parameter space. Thus LMC, and by extension model-averaging, are deeply intertwined with the 516 dynamics of singular values that we explore in Section 4. Full discussion is deferred to Appendix C. 517

518 519

520

7 DISCUSSION

We provide an empirical perspective to understand deep learning through the lens of SVD dynamics. We first note a tendency toward rank minimization on a small scale in grokking, then expand these findings to practical networks and tasks. In addition we find that weight decay, though it explicitly penalizes norm, implicitly promotes this low-rank bias. We also show that generalization and memorization differ in the rank and alignment of solutions found by optimization. We go beyond remarks on generalization and show that magnitude pruning for lottery tickets acts similarly to low-rank pruning, and LMC coincides with the sharing of top singular vectors between checkpoints.

528 While a comprehensive theory for all these results remains elusive, these observations can act as 529 a platform for a deeper understanding of deep learning. Notably, the observed spectral dynamics 530 appear consistent across diverse settings, even without restrictive assumptions like balanced initial-531 ization, linearity, or small weight scales. This suggests a common underlying mechanism.

532 On the empirical side, several interesting problems present themselves. Interpretability of neural 533 networks is a growing area of research (Nanda et al., 2023), and there already exist efforts to interpret 534 singular vectors of convolutional weights (Praggastis et al., 2022). There may also be connections to other unexplained phenomena such as double descent (Belkin et al., 2019; Nakkiran et al., 2021; 535 Davies et al., 2022) or adversarial examples (Szegedy et al., 2013; Ilyas et al., 2019; Hendrycks et al., 536 2021). The solutions to these problems may help design better optimizers or diagnose deployment 537 risks in the wild. There are also concerns of safety (Bai et al., 2022; Mazeika et al., 2024) that better 538 understanding of neural networks can alleviate (Burns et al., 2023; Park et al., 2024). We believe our results contribute another step along this path.

540 REFERENCES

563

565

566

570

577

 Samuel Ainsworth, Jonathan Hayase, and Siddhartha Srinivasa. Git re-basin: Merging models modulo permutation symmetries. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2022.

Maksym Andriushchenko, Francesco D'Angelo, Aditya Varre, and Nicolas Flammarion. Why do we need weight decay in modern deep learning? *arXiv preprint arXiv:2310.04415*, 2023.

Sanjeev Arora, Nadav Cohen, and Elad Hazan. On the optimization of deep networks: Implicit
 acceleration by overparameterization. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2018.

- Sanjeev Arora, Nadav Cohen, Wei Hu, and Yuping Luo. Implicit regularization in deep matrix factorization. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019.
- Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones,
 Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, et al. Constitutional ai: Harm lessness from ai feedback. *arXiv preprint arXiv:2212.08073*, 2022.
- Boaz Barak, Benjamin Edelman, Surbhi Goel, Sham Kakade, Eran Malach, and Cyril Zhang. Hidden progress in deep learning: SGD learns parities near the computational limit. In Advances in Neural Information Processing Systems (NeurIPS), 2022.
- Peter Bartlett. For valid generalization the size of the weights is more important than the size of the network. *Advances in neural information processing systems*, 9, 1996.
 - Mikhail Belkin, Daniel Hsu, Siyuan Ma, and Soumik Mandal. Reconciling modern machinelearning practice and the classical bias–variance trade-off. *Proceedings of the National Academy* of Sciences, 116(32):15849–15854, 2019.
- Emily M Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. On the dangers of stochastic parrots: Can language models be too big? In *Proceedings of the 2021 ACM conference on fairness, accountability, and transparency*, pp. 610–623, 2021.
- Stella Biderman, Hailey Schoelkopf, Quentin Gregory Anthony, Herbie Bradley, Kyle O'Brien, Eric Hallahan, Mohammad Aflah Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward Raff, et al. Pythia: A suite for analyzing large language models across training and scaling. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2023.
- 575 Lukas Biewald. Experiment tracking with weights and biases, 2020. URL https://www.
 576 wandb.com/. Software available from wandb.com.
- Enric Boix-Adserà, Etai Littwin, Emmanuel Abbe, Samy Bengio, and Joshua M Susskind. Trans formers learn through gradual rank increase. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- ⁵⁸¹ Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al. On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*, 2021.
- Johanni Brea, Berfin Simsek, Bernd Illing, and Wulfram Gerstner. Weight-space symmetry in deep
 networks gives rise to permutation saddles, connected by equal-loss valleys across the loss land scape. *arXiv preprint arXiv:1907.02911*, 2019.
- ⁵⁸⁸ Collin Burns, Haotian Ye, Dan Klein, and Jacob Steinhardt. Discovering latent knowledge in language models without supervision. In *The Eleventh International Conference on Learning Representations*, 2023.
- Feng Chen, Daniel Kunin, Atsushi Yamamura, and Surya Ganguli. Stochastic collapse: How gradient noise attracts SGD dynamics towards simpler subnetworks. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2023.

594 Xander Davies, Lauro Langosco, and David Krueger. Unifying grokking and double descent. In 595 NeurIPS ML Safety Workshop, 2022. 596 Sören Dittmer, Emily J King, and Peter Maass. Singular values for ReLU layers. *IEEE Transactions* 597 on Neural Networks and Learning Systems, 31(9):3594–3605, 2019. 598 Simon S Du, Wei Hu, and Jason D Lee. Algorithmic regularization in learning deep homogeneous 600 models: Layers are automatically balanced. In Advances in Neural Information Processing Sys-601 tems (NeurIPS), 2018. 602 Rahim Entezari, Hanie Sedghi, Olga Saukh, and Behnam Neyshabur. The role of permutation 603 invariance in linear mode connectivity of neural networks. In Proceedings of the International 604 Conference on Learning Representations (ICLR), 2021. 605 Tolga Ergen and Mert Pilanci. Path regularization: A convexity and sparsity inducing regularization 607 for parallel relu networks. In Advances in Neural Information Processing Systems (NeurIPS), 608 2023. 609 Simin Fan, Razvan Pascanu, and Martin Jaggi. Deep grokking: Would deep neural networks gener-610 alize better? arXiv preprint arXiv:2405.19454, 2024. 611 612 Ruili Feng, Kecheng Zheng, Yukun Huang, Deli Zhao, Michael Jordan, and Zheng-Jun Zha. Rank 613 diminishing in deep neural networks. In Advances in Neural Information Processing Systems 614 (NeurIPS), 2022. 615 Damien Ferbach, Baptiste Goujaud, Gauthier Gidel, and Aymeric Dieuleveut. Proving linear mode 616 connectivity of neural networks via optimal transport. arXiv preprint arXiv:2310.19103, 2023. 617 618 Stanislav Fort, Gintare Karolina Dziugaite, Mansheej Paul, Sepideh Kharaghani, Daniel M Roy, 619 and Surya Ganguli. Deep learning versus kernel learning: an empirical study of loss landscape 620 geometry and the time evolution of the neural tangent kernel. In Advances in Neural Information 621 Processing Systems (NeurIPS), 2020. 622 Jonathan Frankle and Michael Carbin. The lottery ticket hypothesis: Finding sparse, trainable neural 623 networks. In Proceedings of the International Conference on Learning Representations (ICLR), 624 2018. 625 626 Jonathan Frankle, Gintare Karolina Dziugaite, Daniel Roy, and Michael Carbin. Linear mode connectivity and the lottery ticket hypothesis. In Proceedings of the International Conference on 627 Machine Learning (ICML), 2020. 628 629 Spencer Frei, Gal Vardi, Peter Bartlett, Nathan Srebro, and Wei Hu. Implicit bias in leaky ReLU 630 networks trained on high-dimensional data. In Proceedings of the International Conference on 631 Learning Representations (ICLR), 2022. 632 Tomer Galanti, Zachary S Siegel, Aparna Gupte, and Tomaso Poggio. SGD and weight decay 633 provably induce a low-rank bias in neural networks. arXiv preprint arXiv:2206.05794, 2022. 634 635 Andrey Gromov. Grokking modular arithmetic. arXiv preprint arXiv:2301.02679, 2023. 636 637 Suriya Gunasekar, Blake E Woodworth, Srinadh Bhojanapalli, Behnam Neyshabur, and Nati Srebro. 638 Implicit regularization in matrix factorization. In Advances in Neural Information Processing 639 Systems (NeurIPS), 2017. 640 Stephen Hanson and Lorien Pratt. Comparing biases for minimal network construction with back-641 propagation. In Advances in Neural Information Processing Systems (NeurIPS), 1988. 642 643 Charles R. Harris, K. Jarrod Millman, Stéfan J. van der Walt, Ralf Gommers, Pauli Virtanen, David 644 Cournapeau, Eric Wieser, Julian Taylor, Sebastian Berg, Nathaniel J. Smith, Robert Kern, Matti Picus, Stephan Hoyer, Marten H. van Kerkwijk, Matthew Brett, Allan Haldane, Jaime Fernández 645 del Río, Mark Wiebe, Pearu Peterson, Pierre Gérard-Marchant, Kevin Sheppard, Tyler Reddy, 646 Warren Weckesser, Hameer Abbasi, Christoph Gohlke, and Travis E. Oliphant. Array program-647 ming with NumPy. Nature, 585(7825):357-362, September 2020.

648 Dan Hendrycks, Kevin Zhao, Steven Basart, Jacob Steinhardt, and Dawn Song. Natural adver-649 sarial examples. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern 650 Recognition (CVPR), 2021. 651 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In Advances 652 in Neural Information Processing Systems (NeurIPS), 2020. 653 654 Sepp Hochreiter and Jürgen Schmidhuber. Flat minima. Neural computation, 9(1):1-42, 1997a. 655 656 Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8): 1735–1780, 1997b. 657 658 Minyoung Huh, Hossein Mobahi, Richard Zhang, Brian Cheung, Pulkit Agrawal, and Phillip Isola. 659 The low-rank simplicity bias in deep networks. Transactions on Machine Learning Research, 660 2022. 661 Ahmed Imtiaz Humayun, Randall Balestriero, and Richard Baraniuk. Deep networks always grok 662 and here is why. arXiv preprint arXiv:2402.15555, 2024. 663 J. D. Hunter. Matplotlib: A 2D graphics environment. Computing in Science & Engineering, 9(3): 665 90-95, 2007. 666 Gabriel Ilharco, Marco Tulio Ribeiro, Mitchell Wortsman, Ludwig Schmidt, Hannaneh Hajishirzi, 667 and Ali Farhadi. Editing models with task arithmetic. In Proceedings of the International Con-668 ference on Learning Representations (ICLR), 2022. 669 670 Andrew Ilyas, Shibani Santurkar, Dimitris Tsipras, Logan Engstrom, Brandon Tran, and Aleksander 671 Madry. Adversarial examples are not bugs, they are features. In Advances in Neural Information 672 Processing Systems (NeurIPS), 2019. 673 Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by 674 reducing internal covariate shift. In Proceedings of the International Conference on Machine 675 Learning (ICML), 2015. 676 677 Akira Ito, Masanori Yamada, and Atsutoshi Kumagai. Analysis of linear mode connectivity via 678 permutation-based weight matching. arXiv preprint arXiv:2402.04051, 2024. 679 Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and gener-680 alization in neural networks. In Advances in Neural Information Processing Systems (NeurIPS), 681 2018. 682 683 Ziwei Ji and Matus Telgarsky. Gradient descent aligns the layers of deep linear networks. In Proceedings of the International Conference on Learning Representations (ICLR), 2019. 684 685 Yiding Jiang, Behnam Neyshabur, Hossein Mobahi, Dilip Krishnan, and Samy Bengio. Fantastic 686 generalization measures and where to find them. arXiv preprint arXiv:1912.02178, 2019. 687 688 Keller Jordan, Hanie Sedghi, Olga Saukh, Rahim Entezari, and Behnam Neyshabur. REPAIR: REnormalizing Permuted Activations for Interpolation Repair. In Proceedings of the International 689 Conference on Learning Representations (ICLR), 2022. 690 691 Tero Karras, Miika Aittala, Jaakko Lehtinen, Janne Hellsten, Timo Aila, and Samuli Laine. Analyz-692 ing and improving the training dynamics of diffusion models. arXiv preprint arXiv:2312.02696, 693 2023. 694 Alex Krizhevsky. Learning multiple layers of features from tiny images. Master's thesis, University 695 of Toronto, 2009. 696 697 Anders Krogh and John Hertz. A simple weight decay can improve generalization. In Advances in 698 Neural Information Processing Systems (NeurIPS), 1991. 699 Rohith Kuditipudi, Xiang Wang, Holden Lee, Yi Zhang, Zhiyuan Li, Wei Hu, Rong Ge, and Sanjeev 700 Arora. Explaining landscape connectivity of low-cost solutions for multilayer nets. In Advances in Neural Information Processing Systems (NeurIPS), 2019.

- Tanishq Kumar, Blake Bordelon, Samuel J Gershman, and Cengiz Pehlevan. Grokking as the transition from lazy to rich training dynamics. *arXiv preprint arXiv:2310.06110*, 2023.
- Thien Le and Stefanie Jegelka. Training invariances and the low-rank phenomenon: beyond linear networks. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2021.
- Yann LeCun. The MNIST database of handwritten digits. http://yann.lecun.com/exdb/
 mnist/, 1998.
- Quentin Lhoest, Albert Villanova del Moral, Yacine Jernite, Abhishek Thakur, Patrick von Platen, Suraj Patil, Julien Chaumond, Mariama Drame, Julien Plu, Lewis Tunstall, Joe Davison, Mario Šaško, Gunjan Chhablani, Bhavitvya Malik, Simon Brandeis, Teven Le Scao, Victor Sanh, Canwen Xu, Nicolas Patry, Angelina McMillan-Major, Philipp Schmid, Sylvain Gugger, Clément Delangue, Théo Matussière, Lysandre Debut, Stas Bekman, Pierric Cistac, Thibault Goehringer, Victor Mustar, François Lagunas, Alexander M. Rush, and Thomas Wolf. Datasets: A community library for natural language processing. *arXiv preprint arXiv:2109.02846*, 2021.
- Margaret Li, Suchin Gururangan, Tim Dettmers, Mike Lewis, Tim Althoff, Noah A Smith, and Luke
 Zettlemoyer. Branch-train-merge: Embarrassingly parallel training of expert language models.
 arXiv preprint arXiv:2208.03306, 2022.
- Zhiyuan Li, Yuping Luo, and Kaifeng Lyu. Towards resolving the implicit bias of gradient descent for matrix factorization: Greedy low-rank learning. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2020.
- Ziming Liu, Ouail Kitouni, Niklas S Nolte, Eric Michaud, Max Tegmark, and Mike Williams. Towards understanding grokking: An effective theory of representation learning. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2022.
- Ziming Liu, Eric J Michaud, and Max Tegmark. Omnigrok: Grokking beyond algorithmic data. In
 Proceedings of the International Conference on Learning Representations (ICLR), 2023.
- Ilya Loshchilov and Frank Hutter. Fixing weight decay regularization in Adam. *arXiv preprint arXiv:1711.05101*, 2017.
- 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 *Proceedings of the International Conference on Learning Representations (ICLR)*, 2023.
- Charles H Martin and Michael W Mahoney. Heavy-tailed universality predicts trends in test accuracies for very large pre-trained deep neural networks. In *Proceedings of the SIAM International Conference on Data Mining (ICDM)*, 2020.
- Charles H Martin and Michael W Mahoney. Implicit self-regularization in deep neural networks:
 Evidence from random matrix theory and implications for learning. *Journal of Machine Learning Research*, 22(165):1–73, 2021.
- Mantas Mazeika, Long Phan, Xuwang Yin, Andy Zou, Zifan Wang, Norman Mu, Elham Sakhaee,
 Nathaniel Li, Steven Basart, Bo Li, et al. Harmbench: A standardized evaluation framework for
 automated red teaming and robust refusal. In *Forty-first International Conference on Machine Learning*, 2024.
- 747 Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mixture 748 models. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 749 2016.
- William Merrill, Nikolaos Tsilivis, and Aman Shukla. A tale of two circuits: Grokking as competition of sparse and dense subnetworks. In *ICLR 2023 Workshop on Mathematical and Empirical Understanding of Foundation Models*, 2023.
- Paolo Milanesi, Hachem Kadri, Stéphane Ayache, and Thierry Artières. Implicit regularization in deep tensor factorization. In *Proceedings of the International Joint Conference on Neural Networks (IJCNN)*, 2021.

756 757 758	Sreyas Mohan, Zahra Kadkhodaie, Eero P Simoncelli, and Carlos Fernandez-Granda. Robust and interpretable blind image denoising via bias-free convolutional neural networks. <i>arXiv preprint arXiv:1906.05478</i> , 2019.
759 760 761	Rotem Mulayoff and Tomer Michaeli. Unique properties of flat minima in deep networks. In <i>Proceedings of the International Conference on Machine Learning (ICML)</i> , 2020.
762 763 764	Vaishnavh Nagarajan and J Zico Kolter. Uniform convergence may be unable to explain gener- alization in deep learning. In Advances in Neural Information Processing Systems (NeurIPS), 2019.
765 766 767 768	Preetum Nakkiran, Gal Kaplun, Yamini Bansal, Tristan Yang, Boaz Barak, and Ilya Sutskever. Deep double descent: Where bigger models and more data hurt. <i>Journal of Statistical Mechanics: Theory and Experiment</i> , 2021(12), 2021.
769 770	Neel Nanda, Lawrence Chan, Tom Lieberum, Jess Smith, and Jacob Steinhardt. Progress measures for grokking via mechanistic interpretability. <i>arXiv preprint arXiv:2301.05217</i> , 2023.
771 772 773	Behnam Neyshabur, Ryota Tomioka, and Nathan Srebro. In search of the real inductive bias: On the role of implicit regularization in deep learning. <i>arXiv preprint arXiv:1412.6614</i> , 2014.
774 775	Behnam Neyshabur, Hanie Sedghi, and Chiyuan Zhang. What is being transferred in transfer learn- ing? In Advances in Neural Information Processing Systems (NeurIPS), 2020.
776 777 778	Greg Ongie and Rebecca Willett. The role of linear layers in nonlinear interpolating networks. <i>arXiv</i> preprint arXiv:2202.00856, 2022.
779 780 781	Vassil Panayotov, Guoguo Chen, Daniel Povey, and Sanjeev Khudanpur. Librispeech: an asr corpus based on public domain audio books. In <i>Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)</i> , 2015.
782 783 784	Rahul Parhi and Robert D Nowak. Deep learning meets sparse regularization: A signal processing perspective. <i>IEEE Signal Processing Magazine</i> , 40(6):63–74, 2023.
785 786	Kiho Park, Yo Joong Choe, and Victor Veitch. The linear representation hypothesis and the geometry of large language models. In <i>Forty-first International Conference on Machine Learning</i> , 2024.
787 788 789 790	Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. PyTorch: An imperative style, high-performance deep learning library. In <i>Advances in Neural Information Processing Systems</i> (<i>NeurIPS</i>), 2019.
791 792 793 794	Mansheej Paul, Feng Chen, Brett W Larsen, Jonathan Frankle, Surya Ganguli, and Gintare Karolina Dziugaite. Unmasking the lottery ticket hypothesis: What's encoded in a winning ticket's mask? In <i>Proceedings of the International Conference on Learning Representations (ICLR)</i> , 2022.
795 796 797	Jeffrey Pennington, Samuel Schoenholz, and Surya Ganguli. The emergence of spectral universality in deep networks. In <i>Proceedings of the International Conference on Artificial Intelligence and</i> <i>Statistics (AISTATS)</i> , 2018.
798 799 800 801	Alethea Power, Yuri Burda, Harri Edwards, Igor Babuschkin, and Vedant Misra. Grokking: Gen- eralization beyond overfitting on small algorithmic datasets. <i>arXiv preprint arXiv:2201.02177</i> , 2022.
802 803 804	Brenda Praggastis, Davis Brown, Carlos Ortiz Marrero, Emilie Purvine, Madelyn Shapiro, and Bei Wang. The SVD of convolutional weights: a CNN interpretability framework. <i>arXiv preprint arXiv:2208.06894</i> , 2022.
805 806 807	Xingyu Qu and Samuel Horvath. Rethink model re-basin and the linear mode connectivity. <i>arXiv</i> preprint arXiv:2402.05966, 2024.
808 809	Nasim Rahaman, Aristide Baratin, Devansh Arpit, Felix Draxler, Min Lin, Fred Hamprecht, Yoshua Bengio, and Aaron Courville. On the spectral bias of neural networks. In <i>Proceedings of the International Conference on Machine Learning (ICML)</i> , 2019.

- 810 Prajit Ramachandran, Barret Zoph, and Quoc V Le. Searching for activation functions. arXiv 811 preprint arXiv:1710.05941, 2017. 812 Alexandre Rame, Matthieu Kirchmeyer, Thibaud Rahier, Alain Rakotomamonjy, Patrick Gallinari, 813 and Matthieu Cord. Diverse weight averaging for out-of-distribution generalization. In Advances 814 in Neural Information Processing Systems (NeurIPS), 2022. 815 816 Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-817 conditional image generation with CLIP latents. arXiv preprint arXiv:2204.06125, 2022. 818 Noam Razin and Nadav Cohen. Implicit regularization in deep learning may not be explainable by 819 norms. In Advances in Neural Information Processing Systems (NeurIPS), 2020. 820 821 Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedi-822 cal image segmentation. In Proceeding of the International Conference on Medical Image Com-823 puting and Computer-Assisted Intervention (MICCAI), 2015. 824 Olivier Roy and Martin Vetterli. The effective rank: A measure of effective dimensionality. In 825 Proceedings of the European Signal Processing Conference (EUSIPCO), 2007. 826 827 Ildus Sadrtdinov, Dmitrii Pozdeev, Dmitry P Vetrov, and Ekaterina Lobacheva. To stay or not to 828 stay in the pre-train basin: Insights on ensembling in transfer learning. In Advances in Neural 829 Information Processing Systems (NeurIPS), 2023. 830 A Saxe, J McClelland, and S Ganguli. Exact solutions to the nonlinear dynamics of learning in deep 831 linear neural networks. In Proceedings of the International Conference on Learning Representa-832 tions (ICLR), 2014. 833 Steffen Schotthöfer, Emanuele Zangrando, Jonas Kusch, Gianluca Ceruti, and Francesco Tudisco. 834 Low-rank lottery tickets: Finding efficient low-rank neural networks via matrix differential equa-835 tions. In Advances in Neural Information Processing Systems (NeurIPS), 2022. 836 837 Joseph Shenouda, Rahul Parhi, Kangwook Lee, and Robert D Nowak. Vector-valued variation 838 spaces and width bounds for DNNs: Insights on weight decay regularization. arXiv preprint 839 arXiv:2305.16534, 2023. 840 Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image 841 recognition. arXiv preprint arXiv:1409.1556, 2014. 842 843 Berfin Simsek, François Ged, Arthur Jacot, Francesco Spadaro, Clément Hongler, Wulfram Gerst-844 ner, and Johanni Brea. Geometry of the loss landscape in overparameterized neural networks: 845 Symmetries and invariances. In Proceedings of the International Conference on Machine Learn-846 ing (ICML), 2021. 847 Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised 848 learning using nonequilibrium thermodynamics. In Proceedings of the International Conference 849 on Machine Learning (ICML), 2015. 850 Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, 851 and Rob Fergus. Intriguing properties of neural networks. arXiv preprint arXiv:1312.6199, 2013. 852 853 Zhiquan Tan and Weiran Huang. Understanding grokking through a robustness viewpoint. arXiv 854 preprint arXiv:2311.06597, 2023. 855 Vimal Thilak, Etai Littwin, Shuangfei Zhai, Omid Saremi, Roni Paiss, and Joshua Susskind. The 856 slingshot mechanism: An empirical study of adaptive optimizers and the grokking phenomenon. arXiv preprint arXiv:2206.04817, 2022. 858 859 Nadav Timor, Gal Vardi, and Ohad Shamir. Implicit regularization towards rank minimization in relu networks. In Proceedings of the International Conference on Algorithmic Learning Theory 861 (ALT), 2023. 862 Twan Van Laarhoven. L2 regularization versus batch and weight normalization. arXiv preprint 863
- Twan Van Laarhoven. L2 regularization versus batch and weight normalization. *arXiv preprint arXiv:1706.05350*, 2017.

- Gal Vardi and Ohad Shamir. Implicit regularization in relu networks with the square loss. In *Proceedings of the Annual Conference on Learning Theory (COLT)*, 2021.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In Advances in Neural Information Processing Systems (NeurIPS), 2017.
- 870 Binxu Wang and John Vastola. Stable diffusion, ML from scratch: day 871 2022. URL https://colab.research.google.com/drive/ 2, 872 1Y5wr91g5jmpCDiX-RLfWL1eSBWoSuLqO?usp=sharing#scrollTo= 9is-DXZYwIIi. 873
- Hongyi Wang, Saurabh Agarwal, and Dimitris Papailiopoulos. Pufferfish: Communication-efficient models at no extra cost. *Proceedings of Machine Learning and Systems*, 3:365–386, 2021.
- Zihan Wang and Arthur Jacot. Implicit bias of SGD in l_{2} -regularized linear DNNs: One-way jumps from high to low rank. *arXiv preprint arXiv:2305.16038*, 2023.
- Michael L. Waskom. seaborn: statistical data visualization. Journal of Open Source Software, 6 (60):3021, 2021.
- Mitchell Wortsman, Gabriel Ilharco, Samir Ya Gadre, Rebecca Roelofs, Raphael Gontijo-Lopes, Ari S Morcos, Hongseok Namkoong, Ali Farhadi, Yair Carmon, Simon Kornblith, and Ludwig
 Schmidt. Model soups: averaging weights of multiple fine-tuned models improves accuracy
 without increasing inference time. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2022.
- Zhiwei Xu, Yutong Wang, Spencer Frei, Gal Vardi, and Wei Hu. Benign overfitting and grokking
 in ReLU networks for XOR cluster data. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2023.
- Can Yaras, Peng Wang, Wei Hu, Zhihui Zhu, Laura Balzano, and Qing Qu. Invariant low-dimensional subspaces in gradient descent for learning deep matrix factorizations. In *NeurIPS 2023 Workshop on Mathematics of Modern Machine Learning*, 2023a.
- Can Yaras, Peng Wang, Wei Hu, Zhihui Zhu, Laura Balzano, and Qing Qu. The law of parsimony
 in gradient descent for learning deep linear networks. *arXiv preprint arXiv:2306.01154*, 2023b.
- Hao Yu and Jianxin Wu. Compressing transformers: Features are low-rank, but weights are not! In
 Proceedings of the National Conference on Artificial Intelligence (AAAI), 2023.
- Xiyu Yu, Tongliang Liu, Xinchao Wang, and Dacheng Tao. On compressing deep models by low
 rank and sparse decomposition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- David Yunis, Kumar Kshitij Patel, Pedro Henrique Pamplona Savarese, Gal Vardi, Jonathan Frankle, Matthew Walter, Karen Livescu, and Michael Maire. On convexity and linear mode connectivity in neural networks. In *OPT 2022: Optimization for Machine Learning (NeurIPS 2022 Workshop)*, 2022.

- Emanuele Zangrando, Piero Deidda, Simone Brugiapaglia, Nicola Guglielmi, and Francesco Tud isco. Neural rank collapse: Weight decay and small within-class variability yield low-rank bias.
 arXiv preprint arXiv:2402.03991, 2024.
- Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding
 deep learning (still) requires rethinking generalization. *Communications of the ACM*, 64(3):107–115, 2021.
- Guodong Zhang, Chaoqi Wang, Bowen Xu, and Roger Grosse. Three mechanisms of weight decay regularization. *arXiv preprint arXiv:1810.12281*, 2018a.
- Hongyi Zhang, Yann N Dauphin, and Tengyu Ma. Fixup initialization: Residual learning with out normalization. In *Proceedings of the International Conference on Learning Representations* (*ICLR*), 2018b.

918 919 920	Liu Ziyin, Botao Li, and Xiangming Meng. Exact solutions of a deep linear network. In Advances in Neural Information Processing Systems (NeurIPS), 2022.
921 922	Andy Zou, Zifan Wang, J Zico Kolter, and Matt Fredrikson. Universal and transferable adversarial attacks on aligned language models. <i>arXiv preprint arXiv:2307.15043</i> , 2023.
923	
924	
925	
926	
927	
928	
929	
930	
931	
932	
933	
934	
935	
936	
937	
938	
939	
940	
942	
943	
944	
945	
946	
947	
948	
949	
950	
951	
952	
953	
954	
955	
956	
957	
900	
959	
961	
962	
963	
964	
965	
966	
967	
968	
969	
970	
971	

1.6 0.8 1.4 0.10 0.8 0.98 1.2 0.08 0.6 는 0.6 늡 ≳^{1.0} 0.96 40 aver 0.06 .4 July Nal N 0.8 0.94 0.4 0.04 0.6 0.2 1 97 0.02 0.2 0.4 0.0 0.00 25 Epoch 25 Epoch 50 25 Epoch 50 25 Epoch 25 Epoch 1.00 0.910 2.5 0.8 0.06 0.98 0.905 2.0 0.6 40 0.96 늡 0.900 ≳ 1.5 0.04 Laver .1 U.1 le/ n 94 B 60 0.895 1.0 0.02 0.2 0.890 0.5 ed top err 0.0 14 0.00 50 25 Epoch 50 25 Epoch 25 Epoch 50 25 Epoch Epoch (a) Train Err. (b) Val. Err. (c) SVs (d) Eff. Rank (e) Alignment

Figure 8: Dynamics with random labels for VGG. **Top row:** results with true labels. **Bottom row:** results with random labels. We see that the middle layers have a lower effective rank when using true labels and that alignment in the middle layers persists throughout training. The results are less stark in the VGG case, but similar to the MLP.

A EXPLANATION OF BALANCEDNESS

994 Prior work on deep linear networks (Arora et al., 2019; Milanesi et al., 2021) suggests that rank 995 minimization may describe implicit regularization in deep matrix factorization better than simple 996 matrix norms. See Arora et al. (2018) (Appendix A) for a detailed argument. However, a critical assumption used in these works is "balanced initialization." This means that for consecutive 997 matrices W_i and W_{i+1} in the product matrix $\prod_i W_j$, we have $W_{i+1}^{\top}W_{i+1} = W_i W_i^{\top}$ at initial-998 ization. Decomposing these matrices with SVDs and leveraging orthogonality, this simplifies to 999 $V_{i+1} \Sigma_{i+1}^2 V_{i+1}^{\top} = U_i \Sigma_i^2 U_i^{\top}$ where U_i and V_{i+1} are orthogonal matrices. Since these are orthogonal 1000 decompositions of the same matrix, their diagonals must be equivalent, allowing for the permuta-1001 tion of elements with the same value. This leads to $U_i = V_{i+1}O$ up to signs, where O is a block 1002 diagonal permutation matrix that may permute the rows of equivalent diagonal elements. Notably, if 1003 all diagonal elements are distinct and U_i and V_{i+1} are square matrices, then $U_i = V_{i+1}$ up to signs. 1004 This gives us matching singular vectors for consecutive matrices. 1005

1006 1007

1008

972

973

974

975

976

977

978

979

980

981

982

983

984

985 986

987

988

989

990 991 992

993

B SPECTRAL DYNAMICS WITH RANDOM LABELS

Given the observations connecting generalization and rank thus far, and the enlightening view on the implicit effects of weight decay, we are interested in seeing whether the perspective developed sheds any light on the classic random label memorization experiments of Zhang et al. (2021).

Similar to Zhang et al. (2021), we train a MLP, VGG and an LSTM to fit random or true labels. Please see Appendix D for the details regarding the experimental setup. Zhang et al. (2021) decay the learning rate to zero, and the random label experiments only converge late in training. Consequently, we use a constant learning rate to control this phenomenon. We see in Figure 7 that both cases are able to achieve zero error, though with different singular value evolution and alignment in the middle layer.

1018 Surprisingly in Figure 7, we see that with true labels the inner layers are low rank, while with 1019 random labels they are much higher rank. This may be explained by the shared structure in the 1020 true classes of the dataset, which manifests in the parameters. Even more surprisingly, we find 1021 here that even without weight decay, inner layers align with true labels, while with random labels, this alignment occurs and then disappears with more training. This is particularly intriguing as 1023 there are non-linearities that could theoretically separate the network from the linear case, and yet strong alignment occurs despite that. Such alignment has not yet been leveraged by existing theory, 1024 and might provide structured assumptions for new understanding. Results on the VGG (Figure 8) 1025 are qualitatively quite similar, including on the alignment point. Results on the LSTM (Figure 9)



Figure 9: Dynamics with random labels for LSTM. Top row: results with true labels. Bottom row: 1041 results with random labels. We see that the middle layers have a lower effective rank when using 1042 true labels and that alignment in the middle layers persists throughout training. Though the LSTM 1043 doesn't fit the random labels perfectly, the results are qualitatively similar to the other cases, except 1044 alignment is almost nonexistent.

1047

1049

1051

are weakly similar, though the alignment is much weaker. In summary, these results suggest that viewing generalization through the lens of rank and alignment may be fruitful. 1048

1050 С **BEYOND GENERALIZATION**

1052 We have seen over the course of many experiments that deep models are biased toward low rank, and 1053 that there is a tempting connection between rank minimization and generalization. Still, the lens of 1054 spectral dynamics can be applied more broadly. In the following subsections, we explore two phenomena: lottery tickets (Frankle & Carbin, 2018) and linear mode connectivity (Frankle et al., 2020). 1055 Beyond shedding further light on neural networks, these phenomena have implications for more ef-1056 ficient inference and storage, as well as understanding the importance of pretraining (Neyshabur 1057 et al., 2020). We find that lottery tickets are a sparse approximation of final-checkpoint top singular 1058 vectors. The ability to linearly interpolate between faraway checkpoints and improve performance 1059 coincides strongly with top singular vector sharing between checkpoints. Such observations may form a foundation for a better understanding compression and model averaging (Wortsman et al., 1061 2022; Ilharco et al., 2022).

1062 1063

1064

1069 1070

C.1 TOP SINGULAR VECTORS BECOME STABLE EARLIER

Before we explore the phenomena, we first make another observation that will be helpful. As top singular values grow disproportionately large, it would be natural that top singular vectors become stable in direction as the gradients remain small. To demonstrate this, for a given matrix in the 1067 network $W_i(t) = \sum_{j=1}^R \sigma_j(t) u_j(t) v_j(t)^\top$ at training time t, we compute 1068

$$S(t)_{jk} = |\langle u_j(t)v_j(t)^\top, u_k(T)v_k(T)^\top \rangle|,$$
(5)

1071 where T is the final step of training, and the absolute value is taken to ignore sign flips in the SVD 1072 computation. We then plot the diagonal of this matrix $S(t)_{ii} \forall i \leq 100$ over time. We also use a scalar measure of the diagonal to summarize like in the alignment case: $s(t) = \frac{1}{10} \sum_{i} S(t)_{ii}$. In 1074 Figure 10, we see that top singular vectors converge in direction earlier than bottom vectors.

1075

C.2 LOTTERY TICKETS PRESERVE FINAL TOP SINGULAR VECTORS 1077

As large singular vectors will become stable late in training, we wonder about the connection to 1078 magnitude pruning and the lottery ticket hypothesis. Frankle & Carbin (2018) first showed evidence 1079 for the lottery ticket hypothesis, the idea that there exist sparse subnetworks of neural networks that



Figure 10: Top row: Singular vector agreement for a single matrix in the middle of each model 1099 (diagonal of Eqn. 5). Notice top singular vectors become stable in direction earlier. Bottom row: 1100 Summary score for each matrix across architectures. As we move down the y-axis, the depth of 1101 the parameters in the model increases, while the x-axis tracks training time. The sharp transition 1102 midway through training in the VGG case is likely due to a 10x learning rate decay.

1105 can be trained to a comparable performance as the full network, where the sparse mask is computed 1106 from the largest magnitude weights of the network at the end of training. Frankle et al. (2020) build further on this hypothesis and notice that, for larger networks, the masking cannot begin at 1107 initialization, but rather at some point early in training. Still, the mask must come from the end of 1108 training. 1109

1110 The reason for this particular choice of mask may be connected to the dynamics we previously 1111 observed. Specifically, at the end of training large singular values are disproportionately larger, so 1112 high-magnitude weights may correspond closely to weights in the top singular vectors at the end of training. If magnitude masks were computed at the beginning, the directions that would become 1113 the top singular vectors might be prematurely masked as they have not yet stabilized, which may 1114 prevent learning on the task. 1115

1116 Here we train an unmasked VGG-16 (Simonyan & Zisserman, 2014) on CIFAR10, then compute 1117 either a random mask, or a global magnitude mask from the end of training, and rewind to an early 1118 point (Frankle et al., 2020) to start sparse retraining. We also do the same with an LSTM (Hochreiter & Schmidhuber, 1997b) on LibriSpeech (Panayotov et al., 2015). Please see Appendix D for details. 1119 In Figures 11 and 12, we plot the singular vector agreement (SVA, Eqn. 5) between the final model, 1120 masked and unmasked, where we see exactly that magnitude masks preserve the top singular vectors 1121 of parameters, and with increasing sparsity fewer directions are preserved. Even though prior work 1122 has remarked that it is possible to use low-rank approximations for neural networks (Yu et al., 2017), 1123 and others have explicitly optimized for low-rank lottery tickets (Wang et al., 2021; Schotthöfer 1124 et al., 2022), we rather are pointing out that the magnitude pruning procedure seems to recover a 1125 low-rank approximation. 1126

We also compute the singular vector agreement (SVA) between the masked model trajectory and the 1127 original unmasked model trajectory (diagonal of Eqn. 5). We see in Figures 11 and 12 that there is 1128 no agreement between the bottom singular vectors at all, but there is still loose agreement in the top 1129 singular vectors. Thus, it seems the mask allows the dynamics of only the top singular vectors to 1130 remain similar, which we know are most important from the pruning analysis in Figure 4. 1131

Preserving top singular vectors by pruning seems like a natural outcome of large matrices, so as a 1132 control, we follow exactly the same protocol except we generate the mask randomly with the same 1133 layerwise sparsity. We can see in Figures 11 and 12 that this results in much lower preservation of



Figure 11: Pruning results for VGG. **Top row:** Magnitude pruning. **Bottom row:** random pruning. **First column:** Training loss. We see that at 5% sparsity magnitude pruning is significantly better than random pruning of the same layerwise sparsity. 2nd column: Singular vector alignment pre-and post-pruning at the end of training for a single layer (the 3rd convolution). We see that magnitude pruning approximates the top singular vectors, while random pruning at the same level does not. **3rd column:** Singular vector alignment score pre- and post-pruning across all layers. Agreement is higher across all layers for magnitude pruning, though later layers do not agree, likely as later layers are wider so weights are lower magnitude. 4th column: Singular vector alignment between the pruned and unpruned models along the training trajectory. We see that the magnitude pruning still has similar dynamics in its top singular vectors, while random pruning does not. Last column: Singular vector alignment score between pruned and unpruned models across layers and time. Again evolution is similar for early layers with magnitude pruning, and completely different for random pruning.



Figure 12: Pruning results for LSTM. Top row: Magnitude pruning. Bottom row: random pruning. See Figure 11 for details. Results are quite similar for the LSTM at 25% pruning as the VGG in Figure 11.

top singular vector dynamics, and also performs worse, as in (Frankle et al., 2020). It would not be surprising that random pruning is worse if simply evaluated at the end of training, but masking is applied quite early in training at epoch 4 of 164 long before convergence, so it's striking that the network now fails to learn further even though it is far from convergence. We interpret this as evidence that the mask has somehow cut signal flow between layers, so it is now impossible for the network to learn further, while magnitude pruning and rewinding still allows signals to pass that eventually become important.

1188 C.3 SPECTRAL DYNAMICS AND LINEAR MODE CONNECTIVITY 1189

1190 We come to the final phenomenon that we seek to describe: linear mode connectivity. Linear mode connectivity (LMC) is the property that one can interpolate linearly between two different minima 1191 in weight space and every parameter set along that path performs well, which gives the impression 1192 that the loss surface of neural networks is somehow convex despite its theoretical nonconvexity. 1193 This was first demonstrated in small networks with the same initialization (Nagarajan & Kolter, 1194 2019), then expanded to larger networks and connected to lottery tickets (Frankle et al., 2020; Paul 1195 et al., 2022). Entezari et al. (2021) first conjecture that two arbitrary minima show LMC up to 1196 permutation, and demonstrate it in simple models. This was expanded to wide models (Ainsworth 1197 et al., 2022; Jordan et al., 2022; Qu & Horvath, 2024), and can be proven in various ways (Kuditipudi 1198 et al., 2019; Brea et al., 2019; Simsek et al., 2021; Ferbach et al., 2023), but it does not hold for 1199 standard models (Qu & Horvath, 2024). LMC has also been exploited for model-averaging and performance gains (Wortsman et al., 2022; Ilharco et al., 2022; Rame et al., 2022). Still despite all 1201 of this work, we lack a description for why LMC occurs. In particular: why is there a convex, high dimensional (Yunis et al., 2022) basin that models find shortly in training (Frankle et al., 2020), or 1202 after pretraining (Neyshabur et al., 2020; Sadrtdinov et al., 2023)? We do not answer this question 1203 in full, but find an interesting view through the singular vectors. 1204

LINEAR MODE CONNECTIVITY CORRELATES WITH TOP SINGULAR VECTOR C.3.1 1206 AGREEMENT 1207

1208 As we saw earlier directional convergence of top singular vectors in Figure 10, it suggests the dy-1209 namics of those components are more stable, so we might expect mode-connected solutions to share 1210 these components. To examine this, we plot agreement between the singular vectors of the weight 1211 matrices at either endpoint of branches:

1213
1214
1215
1216
1217

$$W^{(1)}(T) = \sum_{j}^{R} \sigma_{j}(T) u_{j}(T) v_{j}(T)^{\top} ,$$

 $W^{(2)}(T) = \sum_{k}^{R} \sigma_{k}'(T) u_{k}'(T) v_{k}'(T)^{\top} ,$

1214

1212

1205

1215

1216

1217 1218

spawned from the same initialization in training. If the branches are split from an initialization 1219 on a trunk trajectory W(t), we call t the split point or epoch. We visualize the diagonal of 1220 $|\langle u_j(T)v_j(T)^{\top}, u_k'(T)v_k'(T)^{\top}\rangle|_{jk}$ vs. split epoch, where the absolute value is taken to ignore sign 1221 flips in SVD computation. 1222

To remind the reader, LMC only occurs after a small amount of training time has passed. Too early 1223 and the final models of each branch will show a bump, or barrier, in the loss surface along the linear 1224 interpolation (Frankle et al., 2020). To measure this precisely, we use the definition from Neyshabur 1225 et al. (2020), which is the maximum deviation from a linear interpolation in the loss, an empirical 1226 measure for convexity in this linear direction. When this deviation is 0, we consider the checkpoints 1227 to exhibit LMC. Please see Appendix D.10 for details on the calculation. Given evidence in Figure 4 1228 that top components are the most important for prediction, and that top components become stable 1229 before training has finished, it is plausible that LMC is connected to the stability of top singular 1230 vectors in the later portion of training. 1231

This would mean that checkpoints that do not exhibit the LMC property should not share top singular 1232 vectors, while checkpoints that do exhibit the LMC property should share top singular vectors. 1233 We see in Figure 13 that this is the case across models and tasks, where the alignment between 1234 endpoints is much stronger in top singular vectors. We also see no LMC and poor agreement in 1235 top components between branches that have initializations from different trunk trajectories, but with 1236 the same split epoch t and the same branch data order in Figure 14. Thus, these top directions are 1237 not a unique property of the architecture and data, but rather are dependent on initialization. It is notable that concurrent work (Ito et al., 2024) arrives at a similar conclusion: permutation solvers between optima match top singular vectors. Though the conclusions are similar, their experiments 1239 are primarily conducted on smaller scale settings, and only for permutation matching at the end of 1240 training. Here we connect these observations to the optimization behavior of networks throughout 1241 training.

1278

1279 1280



Figure 13: Top row: Barrier size vs. split step. Middle row: singular vector agreement for a single 1269 matrix parameter between branch endpoints that share a common trunk. Bottom row: summary 1270 statistic for singular vector agreement across layers vs. split step. We see that as models exhibit LMC, they also share top singular vectors.

PERTURBING BREAKS LINEAR MODE CONNECTIVITY AND SINGULAR VECTOR C.3.2 AGREEMENT SIMULTANEOUSLY

1281 To make the connection between top singular vectors and LMC even tighter, we intervene in the 1282 normal training process. If we add random perturbations to destabilize the components that will 1283 become the top components long before they have converged, and if singular vector agreement is 1284 tied to LMC, we would like to see that final models no longer exhibit the LMC property. Indeed 1285 this is the case. In Figure 15, when increasingly large random perturbations are applied, the barrier 1286 between final checkpoints increases and the LMC behavior disappears. Please see Appendix D 1287 for details. In addition, the previously-strong singular vector agreement disappears simultaneously. Thus it seems this agreement is tied to linear mode connectivity. 1288

1289 We speculate that, due to the results in Figure 4 that show the top half of the SVDs are much 1290 more critical for performance, if these components are shared then interpolating will not affect 1291 performance much. Rather, interpolation will eliminate the orthogonal bottom components which 1292 may only make a minor impact on performance. If however the top components are not shared, 1293 then interpolating between two models will remove these components, leading to poor performance in between. Such observations may help in explaining the utility of pretraining (Neyshabur et al., 1294 2020), weight averaging (Rame et al., 2022; Wortsman et al., 2022; Ilharco et al., 2022) or the use 1295 of LoRA (Huh et al., 2022) to replace full finetuning.



Figure 14: **Top row:** Barrier size vs. split step. **Middle row:** singular vector agreement for a single matrix parameter between branch endpoints that do not share a common trunk, but do share split time and branch data order. **Bottom row:** summary statistic for singular vector agreement across layers. We see that when branches do not share a common trunk, there is neither LMC nor singular vector agreement, even though the optimization is otherwise the same.

1330

1329 D EXPERIMENTAL DETAILS

For all experiments, we use 3 random seeds and average all plots over those 3. This is relatively small, but error bars tend to be very tight, and due to the high volume of runs required for this work we lack the resources to run much more.

In order to compute alignment we consider only pairs of layers that directly feed into each other, and ignore the influence of residual connections so as to cut down on the number of comparisons.
 Specifics on individual architectures are given below.

1337

1338 D.1 IMAGE CLASSIFICATION WITH VGG

We train a VGG-16 (Simonyan & Zisserman, 2014) on CIFAR-10 (Krizhevsky, 2009) for 164 epochs, following hyperparameters and learning rate schedule in (Frankle et al., 2020), but without data augmentation. For the optimizer we use SGD with batch size 128, initial learning rate 0.1 and momentum of 0.9. We also decay the learning rate 3 times by a factor of 10 at epoch 82, epoch 120, and finally at epoch 160. We also use a minor amount of weight decay with coefficient 0.0001.

1345 VGG-16 uses ReLU activations and batch normalization (loffe & Szegedy, 2015), and includes both 1346 convolutional and linear layers. For linear layers we simply compute the SVD of the weight matrix. 1347 For convolutional layers, the parameters are typically stored as a 4D tensor of shape (c_{out}, c_{in}, h, w) 1348 for the output channels, input channels, height and width of the filters respectively. As the filters 1349 compute a transformation from each position and input channel to an output channel, we compute 1349 the SVD of the flattened tensor $(c_{out}, c_{in} \cdot h \cdot w)$, which maps all inputs to outputs, similar to Praggastis



Figure 15: **Top row:** Barrier size vs. perturbation magnitude. **Middle row:** singular vector agreement for a single matrix parameter between branch endpoints vs. perturbation magnitude. **Bottom row:** summary statistic for singular vector agreement across layers with perturbation magnitude. We see that whereas without perturbation models would exhibit LMC after training, with increasing perturbations the LMC property disappears simultaneously with the agreement in top singular vectors.

et al. (2022). This is not the SVD of the entire transformation of the feature map to the next feature
map, but rather the transformation from a set of adjacent positions to a particular position in the next
layer. For the individual SV evolution plot, we use the 12th convolutional layer.

1387 In order to compute alignment of bases between consecutive convolutional layers, $V_{i+1}^{\perp}U_i$ we need 1388 to match the dimensionality between U_i and V_{i+1} . For convolutional layers we are presented with 1389 a question as to how to handle the spatial dimensions h and w as naively the input dimension of 1390 the next layer will be a factor of $h \cdot w$ larger dimension. We experimented with multiple cases, 1391 including aligning at each spatial position individually or averaging over the alignment at all spatial positions, and eventually settled at aligning the output of one layer to the center spatial input of the 1392 next layer. That is, for a 3x3 convolution mapping to a following 3x3 convolution, we compute the 1393 alignment only for position (1,1) of the next layer. This seemed reasonable to us as on average the 1394 edges of the filters showed poorer alignment overall. For the individual alignment plot, we use the 1395 alignment between the 11th and 12th convolutional layers at the center spatial position of the 12th 1396 convolutional layer. 1397

1398

1399 D.2 IMAGE GENERATION WITH UNETS

1400

We train a UNet (Ronneberger et al., 2015) diffusion model (Sohl-Dickstein et al., 2015; Ho et al., 2020) on MNIST (LeCun, 1998) generation. We take model design and hyperparameters from (Wang & Vastola, 2022). In particular we use a 4-layer residual UNet and train with AdamW (Loshchilov & Hutter, 2017) with batch size 128, and learning rate of 0.0003 for 100

26

epochs. This model uses swish (Ramachandran et al., 2017) activations and a combination of linear and convolutional, as well as transposed convolutional layers.

Computing SVDs and alignment is similar to the image classification case described above, except in the case of the transposed convolutions where an extra transpose of dimensions is needed as parameters are stored with the shape (c_{in}, c_{out}, h, w) . For the individual SV evolution plot, we use the 3rd convolutional layer. For the alignment plot, we use the alignment between the 3rd and 4th convolutional layers at the center spatial position of the 4th convolutional layer.

1411 1412

1412 D.3 SPEECH RECOGNITION WITH LSTMs

We train a bidirectional LSTM (Hochreiter & Schmidhuber, 1997a) for automatic speech recognition
on LibriSpeech (Panayotov et al., 2015). We tune for a simple and well-performing hyperparameter
setting. We use AdamW (Loshchilov & Hutter, 2017) with batch size 32, learning rate 0.0003 and
weight decay 0.1 for 50 epochs. We also use a cosine annealing learning rate schedule from 1 to 0
over the entire 50 epochs.

1419 The LSTM only has matrix parameters and biases, so it is straightforward to compute SVDs of 1420 the matrices. For individual SV evolution plots, we plot the 3rd layer input parameter. In the case 1421 of alignment, we make a number of connections: first down depth for the input parameters, then 1422 connecting the previous input parameter to the current hidden parameter in both directions, then 1423 connecting the previous hidden parameter to the current input parameter. In particular the LSTM parameters are stored as a stack of 4 matrices in PyTorch, and we find alignment is highest for 1424 the "gate" submatrix, so we choose that for all plots. For the individual layer alignment, we plot 1425 alignment between the 3rd and 4th layer input parameters. 1426

- 1427
- 1428 D.4 LANGUAGE MODELING WITH TRANSFORMERS

We train a Transformer (Vaswani et al., 2017) language model on Wikitext-103 (Merity et al., 2016). We base hyperparameter choices on the Pythia suite (Biderman et al., 2023), specifically the 160 million parameter configuration with sinusoidal position embeddings, 12 layers, model dimension 768, 12 attention heads per layer, and hidden dimension 768. We use AdamW (Loshchilov & Hutter, 2017) with batch size 256, learning rate 0.0006 and weight decay 0.1. We use a context length of 2048 and clip gradients to a maximum norm of 1. We also use a learning rate schedule with a linear warmup and cosine decay to 10% of the learning rate, like Biderman et al. (2023).

For SVDs, for simplicity we take the SVD of the entire $(3d_{model}, d_{model})$ parameter that computes queries, keys and values from the hidden dimension inside the attention layer, without splitting into individual heads. This is reasonable as the splitting is done after the fact internally. We also take the SVD of the output parameters, and linear layers of the MLPs, which are 2 dimensional matrices. For the individual SV evolution plot, we plot the SVs of W_1 of the 8th layer MLP

For alignment, we consider the alignment of W_Q and W_K matrices, W_V and W_O matrices, computing alignment between heads individually then averaging over all heads. We also consider the alignment between W_O and W_1 of the MLP block, between W_1 and W_2 of the MLP block, and between W_2 and the next attention layer. For the individual layer alignment, we plot alignment between W_1 and W_2 of the 8th layer MLP.

1446 1447

1448

D.5 SPECTRAL DYNAMICS WITH SCALE (PYTHIA)

Here we apply the perspective developed in Section 4 to larger scale models. As we lack the resources to train these models ourselves, we leverage the Pythia (Biderman et al., 2023) family which provides training trajectories for language models across a range of scales (70m to 12b parameters). We are further constrained to the 2.8b parameter model at the largest due to memory requirements when computing SVDs and alignment.

In Figure 16, we see similar rank dynamics across a variety of scales. We choose to select the 7th layer MLP to compare between models as it is present at all scales. We do see an unequal evolution in singular values, but also a contraction as training proceeds for longer. The difference between scales is not very obvious, but slightly fewer of the singular values evolve to be large in the 2.8b model as opposed to the 410m model, which one can see from the thickness of the light magenta



Figure 16: Spectral dynamics of Pythia suite. From top to bottom we examine the 160m, 410m, 1487 1.4b and 2.8b parameter models. Notably, much less noise appears in the alignment plot with increasing scale. Presumably this could be due to the fact that larger dimensional vectors have higher probability to be orthogonal, which may play a role in making optimization easier. We see stronger alignment score (Eqn. 4) in all layers in the larger model, perhaps because of that cleaner signal.

1493

1494 1495

1497

color. The lack of alignment except for the top rank is quite consistent with earlier observations, and such alignment happens much later for the largest model.

1496 D.6 WEIGHT DECAY EXPERIMENTS

All tasks are trained in exactly the same fashion as mentioned previously, with increasing weight decay in the set {0,0.0001,0.001,0.01,0.1,1.0,10.0}. For ease of presentation we consider a subset of settings across tasks. In Figure 18 we include trained model performance and pruned model performance to show that, even with high levels of weight decay, models do not entirely break down. More so, the approximation of the pruned model to the full model gets better with higher weight decay.

1503 1504

1505 D.7 GROKKING EXPERIMENTS

For the Trasnformer, we mostly follow the settings and architecture of Nanda et al. (2023), except we use sinusoidal positional encodings instead of learned.

- For the slingshot case we follow hyperparameter settings in Thilak et al. (2022), Appendix B except with the 1-layer architecture from Nanda et al. (2023) instead of the 2-layer architecture specified.
 W perform addition modulo 97. The original grokking plot in Thilak et al. (2022) appears much more dramatic as it log-scales the x-axis, which we do not do here for clarity.
 - 28



Figure 17: Diagonal of alignment for a single pair over time (Eqn. 3) and alignment metric across 1531 pairs of matrices over time (Eqn. 4) where the y-axis represents depth. From top to bottom, for VGG 1532 we use coefficients $\{0, 0.001, 0.01, 0.1\}$, while for other networks we use coefficients $\{0, 0.1, 1, 10\}$. 1533 We see that the maximum alignment magnitude is higher with large weight decay, and in particular, 1534 the Transformer has the strongest alignment even when nonlinearities separate the MLP layers.

1535 1536

1537 In the case of the deep MLP, we follow Fan et al. (2024), where we use a 12-layer MLP with ReLU 1538 activations and width 400, trained on MSE loss on MNIST (LeCun, 1998). We use 2000 examples, a batch size of 100, weight decay 0.01, and initialization scale 8 (Liu et al., 2023). 1539

1540

1541 D.8 **RANDOM LABEL EXPERIMENTS** 1542

1543 We train a 4-layer MLP on CIFAR10 (Krizhevsky, 2009) with either completely random labels, or the true labels. We use SGD with momentum of 0.9 and constant learning rate of 0.001, and train 1544 for 300 epochs to see the entire trend of training. The major difference to the setting of Zhang et al. 1545 (2021) is the use of a constant learning rate, as their use of a learning rate schedule might conflate 1546 the results. 1547

1548 For the VGG case, we follow our previous hyperparameters, except we leave out weight decay and 1549 learning rate scheduling, instead using a constant learning rate of 0.01.

1550 For the LSTM case, we follow our previous hyperparameters, and extend the training budget to 200 1551 epochs allow for the random label setting to train longer. In this case, our network does not have 1552 sufficient capacity to memorize the data completely. 1553

- 1554 1555
- D.9 MAGNITUDE PRUNING EXPERIMENTS

1556 We use the same VGG setup as described previously. In this case we train til the end, then compute a 1557 global magnitude mask. To do this we flatten all linear and convolutional weights into a single vector, 1558 except for the last linear layer, and sort by magnitude. Then we keep the top 5% of weights globally, 1559 and reshape back to the layerwise masks. This results in different sparsity levels for different layers, 1560 so when generating the random masks, we use the per-layer sparsities that resulted from the global 1561 magnitude mask.

- To retrain the network, we rewind to epoch 4, then continue training with the mask, always setting 1563 other weights and their gradients to 0. We average all results over 3 random seeds. 1564
- For the LSTM we follow exactly the same procedure, except our mask only reaches a level of 25% 1565 sparsity, due to large performance degradations past that.



Figure 18: Training loss over time, where the rows use differing amounts of weight decay. From top to bottom, for VGG we use coefficients $\{0, 0.001, 0.01, 0.1\}$, while for other networks we use coefficients $\{0, 0.1, 1, 10\}$. We see that it is still possible to achieve low training loss under high weight decay, and as we increase the amount of weight decay, the gap between pruned and unpruned parameters closes, lending support to the idea that the parameters become lower rank.

D7 D.10 LMC EXPERIMENTS

1606

1614

1615

We save 5 evenly-spaced checkpoints in the first epoch, as well as at the end of the next 4 epochs for 10 initializations in total. We train 3 trunks, and split 3 branches from each trunk for a total of 9 branches which we average all plots over.

Following Neyshabur et al. (2020), we compute the barrier between checkpoints as follows: given $W^{(1)}(T)$ and $W^{(2)}(T)$ that were branched from W(t) we compute

$$b(t) = (\max_{\alpha \in [0,1]} \mathcal{L}((1-\alpha)W^{(1)}(T) + \alpha W^{(2)}(T)) - ((1-\alpha)\mathcal{L}(W^{(1)}(T)) + \alpha \mathcal{L}(W^{(2)}(T)))$$
(6)

when this quantity is 0, we consider the checkpoints to exhibit LMC.

We recompute batch normalization parameters after interpolating for VGG-16, and group normal ization parameters for the UNet, as these do not necessarily interpolate well (Frankle et al., 2020).
 We also compute singular vector agreement for the same parameter between either branch endpoint.



Figure 19: Grokking and Spectral Dynamics in Modular addition. Top row: 30% data and no 1655 weight decay. 2nd row: 30% data and weight decay 1.0 (grokking), using hyperparameters from 1656 Nanda et al. (2023). 3rd row: 70% data with no weight decay (slingshot), using hyperparameters 1657 from Thilak et al. (2022). Bottom row: 90% data and no weight decay. 1st column: Training and 1658 validation error. 2nd column: Singular value evolution is visualized for the first attention parameter, 1659 where each line represents a single singular value and the color represents the rank. **3rd column:** Effective rank of all layers (Eqn. 1). 4th column: Alignment (Eqn. 3) between the embedding 1660 and the first attention parameter is also visualized, where the y-axis corresponds to index i of the 1661 diagonal. One can see that grokking co-occurs with low-rank weights. In addition, there is an 1662 alignment that begins early in training that evolves up the diagonal. Without weight decay and with 1663 less data, neither grokking nor the other phenomena occur during the entire training budget, but using 1664 more data, even without weight decay, leads to low-rank solutions from the beginning of training. 1665 The slingshot case follows a similar trend, though the validation loss is gradually fit. Across cases 1666 with good generalization, parameters are lower rank, and alignment is also more prevalent in the top ranks.

- 1668
- 1669
- 1670
- 1671

To plot the singular vector (dis)agreement and LMC between different modes, we make 11 evenly
 spaced measurements interpolating between branch endpoints that had the same split epoch, and the same branch seed, but different trunk initializations.

1674 D.11 PERTURBED LMC EXPERIMENTS

1676 We perturb all weights W after the point of dynamics stability where we expect to see LMC at the 1677 end of training (epoch 4 is sufficiently late in all cases) using randomly sampled normal perturbations 1678 $\epsilon \sim \mathcal{N}(0, I)$ with $\|\epsilon\| = \eta \|W\|$ where $\eta \in \{0.0, 0.1, 0.25, 0.5, 1.0, 2.5\}$. We do not perturb the 1679 output layer, as this has a very substantial effect on the optimization. We also do not perturb the 1680 input layer for the Transformer as it is too computationally expensive for our resources.

1681

1682 E LIMITATIONS

1683

There are a few key limitations to our study. As mentioned, we lack the computational resources to run more than 3 random seeds per experiment, though we do find error bars to be quite tight in general (except for the generalization epoch in the grokking experiments). In addition, as discussed we ignore 1D parameters in the neural networks, which may be particularly crucial (especially normalization). In addition, due to computational constraints we do not consider alignment of layers across residual connections as this quickly becomes combinatorial in depth, thus there may be other interesting interactions that we do not observe. Finally, due to computational constraints we are unable to investigate results on larger models than the 12 layer Transformer, which may have different behavior.

1692

1695

1694 F COMPUTE RESOURCES

All experiments are performed on an internal cluster with on the order of 100 NVIDIA 2080ti GPUs or newer. All experiments run on a single GPU in less than 8 hours, though it is extremely helpful to parallelize across machines. We estimate that end-to-end it might take a few days on these resources to rerun all of the experiments in this paper. Additionally, the storage requirements for all of the checkpoints will take on the order of 5 terabytes.

1701

1702 G CODE SOURCES

1703

We use PyTorch (Paszke et al., 2019) and NumPy (Harris et al., 2020) for all experiments and
Weights & Biases (Biewald, 2020) for experiment tracking. We make plots with Matplotlib (Hunter,
2007) and Seaborn (Waskom, 2021). We also use HuggingFace Datasets (Lhoest et al., 2021) for
Wikitext-103 (Merity et al., 2016).

1712 1713

- 1720
- 1722
- 1723
- 1724
- 1725

1726 1727