OPTIMIZER CHOICE MATTERS FOR THE EMERGENCE OF NEURAL COLLAPSE

Anonymous authorsPaper under double-blind review

000

001

002003004

010 011

012

013

014

015

016

018

019

021

024

025

026

027

028

031

033

034

037

038

040 041

042

043

044

046

047

048

051

052

ABSTRACT

Neural Collapse (NC) refers to the emergence of highly symmetric geometric structures in the representations of deep neural networks during the terminal phase of training. Despite its prevalence, the theoretical understanding of NC remains limited. Existing analyses largely ignore the role of the optimizer, thereby suggesting that NC is universal across optimization methods. In this work, we challenge this assumption and demonstrate that the choice of optimizer plays a critical role in the emergence of NC. The phenomenon is typically quantified through NC metrics, which, however, are difficult to track and analyze theoretically. To overcome this limitation, we introduce a novel diagnostic metric, NC0, whose convergence to zero is a necessary condition for NC. Using NC0, we provide theoretical evidence that NC cannot emerge under decoupled weight decay, as implemented in AdamW. Concretely, we prove that SGD, SignGD with coupled weight decay (a special case of Adam), and SignGD with decoupled weight decay (a special case of AdamW) exhibit qualitatively different NC0 dynamics. Finally, we conduct extensive empirical experiments consisting of 3,900 training runs across various datasets, architectures, optimizers, and hyperparameters, confirming our theoretical results. This work provides the first theoretical explanation for optimizer-dependent emergence of NC and highlights the overlooked role of weight-decay coupling in shaping the implicit biases of optimizers.

1 Introduction

Neural networks have driven many of the recent breakthroughs in artificial intelligence, yet the mechanisms underlying their success remain only partially understood. A key empirical clue is neural collapse (NC) – first documented by Papyan et al. (2020) – in which the last-layer feature vectors and classifier weights self-organise into a highly symmetric configuration during the terminal phase of training (TPT). While the reasons for the emergence of NC are still not fully understood, its impact on the behavior of a model is evident. For instance, Liu et al. (2023) induce NC to improve generalization in class-imbalanced training and Galanti et al. (2021) show that the emergence of NC improves transfer learning as well. Furthermore, the presence of NC has been connected to better out-of-distribution detection (Liu & Qin, 2023).

Theoretical explanations for NC have primarily relied on simplified models and assumptions (Mixon et al., 2022; Zhu et al., 2021) that have largely ignored the role of the optimizer, thereby suggesting that NC is universal across optimization methods. In this work, we challenge this assumption and demonstrate that the choice of optimizer plays a critical role in the emergence of NC. Concretely, we show that training with AdamW (Loshchilov & Hutter, 2019) does not lead to an NC solution, whereas training with SGD or Adam (Kingma & Ba, 2014) does. Through extensive experiments, we trace this back to how weight decay is applied in both optimizer and identify the coupling of weight decay as a necessity for the emergence of NC.

One major challenge in studying NC lies in the original metrics, which are difficult to track and analyze theoretically. These metrics were designed to quantify the progressive geometric alignment associated with NC and are expected to converge to zero in the idealized setting where NC holds as training time approaches infinity. However, under realistic training regimes, such as finite training epochs and learning rate decay, these metrics typically plateau at small but nonzero values. As a result, there is no rigorous criterion for determining whether NC has truly occurred.

This limitation motivates us to introduce a novel diagnostic metric, NC0, whose convergence to zero is necessary (though not sufficient) for NC. Unlike previous metrics, NC0 enables a more definitive assessment: if NC0 diverges during training, we can conclude that NC can not occur—even in cases where other NC metrics misleadingly converge to small positive values, creating an illusion of collapse. We discuss the peculiarity of interpreting NC metrics in practice later in Section 4.1. Furthermore, NC0 allows us to go beyond loss landscape analysis and theoretically derive convergence rates with which NC0 converges to zero.

Contribution In this paper, we conduct extensive experiments – spanning over 3,900 training runs – to investigate the role of coupled weight decay in the emergence of NC. We identify coupled weight decay as a key driver of NC in realistic settings, extending recent theoretical insights (Pan & Cao, 2024; Jacot et al., 2024) that were limited to quasi-optimal solutions in simplified models. In particular, we show that the form of weight decay used in adaptive optimizers such as Adam (Kingma & Ba, 2014) and AdamW (Loshchilov & Hutter, 2019) critically affects whether NC emerges. Strikingly, while networks trained with Adam often exhibit NC, AdamW – despite its algorithmic similarity –fails to produce NC, with the corresponding metrics failing to converge to zero over time (Figure 1). This subtle yet consequential distinction has been largely overlooked in prior work.

In summary, we make the following contributions:

- 1. Across a wide range of experiments, we find that coupled weight decay is a necessary condition for NC to emerge.
- 2. Furthermore, we show the accelerating effect of momentum on NC when trained with SGD, being the first result concerning momentum in the context of NC.
- 3. We support our empirical findings with the following theoretical statements on the new NC0 metric:
 - with SGD, NC0 converges to zero at an exponential rate proportional to the weight decay;
 - with sign gradient descent (SignGD) with decoupled weight decay, a special case of AdamW, NC0 converges to some positive constant;
 - with SignGD with coupled weight decay, a special case of Adam, NC0 exhibits a non-monotonic trajectory, increasing before eventually decreasing.

Organization This paper is organized as follows. In Section 2, we recapitulate the four properties to characterize NC and introduce a novel NC property NC0. In Section 3 we present our main experimental results with theoretical support. Finally, Section 4 provides insights and discussions on the implications of our results.

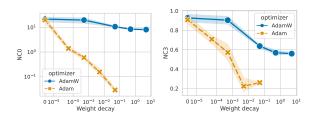


Figure 1: NC0 (left) and NC3 (right) metrics at the end of training. Lower values indicate stronger NC. AdamW shows consistently higher metrics than Adam. Averages computed over runs with varying learning rates and momentum; shaded regions show ± 1 standard deviation. X-axis is log-scaled.

Notation We use $[K] = \{1, 2, ..., K\}$ to denote the index set for any integer $K \in \mathbb{N}$. For a matrix \mathbf{W} , we let $\mathrm{Vec}(\mathbf{W})$ denote the vectorization of \mathbf{W} obtained by stacking its columns. The Frobenius inner product between two matrices \mathbf{W}, \mathbf{W}' is denoted by $\langle \mathbf{W}, \mathbf{W}' \rangle = \mathrm{Tr}(\mathbf{W}^{\top}\mathbf{W}')$. With slight abuse of notation, we write $\|\mathbf{W}\| = \|\mathbf{W}\|_F$ for the Frobenius norm when \mathbf{W} is a matrix, and $\|\mathbf{v}\| = \|\mathbf{v}\|_2$ for the Euclidean norm when \mathbf{v} is a vector. In other words, $\|\mathbf{W}\| = \|\mathrm{Vec}(\mathbf{W})\|$. We denote by \mathbf{I} the identity matrix, by $\mathbf{1}$ the all-ones column vector, and by \mathbf{J} the all-ones matrix, i.e., $\mathbf{J} = \mathbf{1}\mathbf{1}^{\top}$.

2 NEURAL COLLAPSE

Neural collapse (NC), observed during the terminal phase of training (TPT) in deep neural networks (DNN), manifests itself through several geometric properties involving the last-layer features and weights in the *K*-class classification task:

$$\min_{\mathbf{W},\theta} \sum_{n=1}^{N} \ell(\mathbf{W} h_{\theta}(\mathbf{x}_n), y_n) + \frac{\lambda}{2} ||\mathbf{W}||^2 + \frac{\lambda}{2} ||\mathbf{Vec}(\theta)||^2$$
 (1)

where $(\mathbf{x}_n, y_n)_{n=1}^N \subset \mathbb{R}^D \times [K]$ is the training set, $\mathbf{W} \in \mathbb{R}^{K \times P}$ is the last-layer weights, $h_{\theta}(\mathbf{x}_n) \in \mathbb{R}^P$ is the last-layer feature as the output of some backbone parameterized by $\theta, \ell : \mathbb{R}^K \times [K] \to [0, \infty)$ is the loss function, and $\lambda > 0$ is the L2-regularization constant.

These properties, formalized by their corresponding metrics in the original paper Papyan et al. (2020), are:

- NC1 Variability Collapse: Features collapse to their respective class means, indicating that within-class variability vanishes.
- 2. **NC2 Convergence of Class Means to Simplex ETF:** Class means converge to a simplex equiangular tight frame (ETF).
- 3. NC3 Convergence to Self-Duality: Rows of the last-layer weight $\mathbf{W} \in \mathbb{R}^{K \times P}$ align with the columns of the class means, creating a dual relationship between weights and features.
- 4. NC4 Simplification to Nearest-Class-Center: The classifier's decision boundaries are simplified to those of a nearest-class-mean (NCC) classifier.

A solution satisfying all of these properties is referred to as a *NC solution*. In addition to these prior NC properties, we introduce another novel NC property **NC0**, whose convergence to zero is a necessary condition (though not sufficient) for NC.

NC0 - Zero Row Sum of Last-Layer Weight: The row sum of the last-layer weight W in the model converges to zero.

The first observation is that NC0 is a necessary condition for NC2 and NC3:

Proposition 2.1. *NC2 and NC3 implies NC0.*

Proof. For each class $k \in [K]$, we define the class mean $\mu_k = \frac{1}{|\{n:y_n=k\}|} \sum_{n:y_n=k} h_{\theta}(\mathbf{x}_n) \in \mathbb{R}^P$ and the centered class mean $\bar{\mu}_k = \mu_k - \frac{1}{N} \sum_{n=1}^N h_{\theta}(\mathbf{x}_n)$. We concatenate them into a matrix $\mathbf{M} = (\bar{\mu}_k)_{k=1}^K \in \mathbb{R}^{P \times K}$ with $\mathbf{M}\mathbf{1} = 0$, since we centered the class means. By NC2, \mathbf{M} converge to a simplex ETF in the ambient space \mathbb{R}^P , meaning $\mathbf{M}/\|\mathbf{M}\|_F \to \mathbf{Q}\mathbf{M}^*$ where $\mathbf{M}^* \in \mathbb{R}^{K \times K}$ is a unit matrix with columns forming a K-simplex EFT in \mathbb{R}^K and $\mathbf{Q} \in \mathbb{R}^{P \times K}$ is the isometric injection map into the ambient space. Since $\mathbf{M}\mathbf{1} = \mathbf{0}$ and \mathbf{Q} is injective, the unit matrix \mathbf{M}^* has to be in the form: $\mathbf{M}^* \stackrel{\text{def.}}{=} \frac{1}{\sqrt{K-1}} \left(\mathbf{I} - \frac{1}{K}\mathbf{J}\right)$. Hence

$$\mathbf{M}^{\top}\mathbf{M}/\|\mathbf{M}^{\top}\mathbf{M}\|_F^2 \to (\mathbf{Q}\mathbf{M}^*)^{\top}\mathbf{Q}\mathbf{M}^* = (\mathbf{M}^*)^2 = \mathbf{M}^*.$$

On the other hand, NC3 states that $\mathbf{M}/\|\mathbf{M}\| - \mathbf{W}^\top/\|\mathbf{W}\| \to 0$ as $t \to \infty$. Hence we have $\frac{\mathbf{W}\mathbf{W}^\top}{\|\mathbf{W}\|_F^2} - \mathbf{M}^* \to 0$ as $t \to \infty$. Now note that $\mathbf{1}^\top \mathbf{M}^* \mathbf{1} = 0$, hence $\|\mathbf{W}^\top \mathbf{1}\|^2 = \mathbf{1}^\top \mathbf{W} \mathbf{W}^\top \mathbf{1} \to 0$. Note that the last line holds if and only if NC0 holds.

NC0 offers two key advantages. First, it serves as a diagnostic tool: if NC0 does not converge, then at least one of NC2 or NC3 must fail, providing a clear signal that neural collapse cannot occur. Second, NC0 is more mathematically tractable than the original NC metrics, whose dynamics are difficult to analyze and remain underexplored. As we demonstrate in Section 3, NC0's evolution during training can be reliably tracked and used to explain empirical trends observed across different optimizers. In addition, our extensive experiments also show that NC0 is correlating well with prior NC metrics, particularly for small learning rates (see Figure 2). For a more detailed explanation and formal definitions of NC properties and their metrics, we refer the reader to Section B.

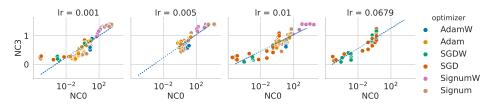


Figure 2: NC0 weakly correlates with NC3 across different optimizers and learning rates. Details on the regression fit can be found in Appendix D.3

3 MAIN RESULT

3.1 EXPERIMENTAL SETUP

We conducted extensive experiments training a ResNet9 and VGG9 using various optimizers, including Adam, AdamW, SGD, SGD with decoupled weight decay (SGDW), Signum (Bernstein et al., 2018), and Signum with decoupled weight decay (SignumW) trained on MNIST, FashionMNIST and Cifar10. Every optimizer is trained with three different learning rates (LR), six different values of momentum, and six different values of weight decay to also control the effect of hyperparameters on the emergence of NC. This resulted in a total of $2 \times 3 \times 6 \times 108 = 3$, 888 training runs. All networks were trained for 200 epochs using a batch size of 128, with the learning rate being decayed by a factor of 10 after one-third and two-thirds of the training duration, as described in the original work by Papyan et al. (2020). In addition, we conducted ablation studies to control for the number of training epochs and to verify that the results also hold for unconstrained feature models (UFM)¹, leading to a total of over 3,900+ training runs. Further details and all experimental results can be found in Appendix D. Ablation studies on the effect of training epochs can be found in Appendix D.4.1

Table 1: Final NC metrics for the same setting as in Figure 5, following the setup of Papyan et al. (2020). Lower values (\downarrow) indicate stronger neural collapse. Values in parentheses represent percentages relative to the metric at initialization.

Optimizer	$\mathrm{NC0}_{\downarrow}$	NC1↓	NC2↓	NC3↓
SGD	2.14e-04 (< -99.5%)	0.05 (-99.3%)	0.29 (-63.0%)	0.35 (-75.1%)
SGDW	0.55 (-68.9%)	0.26 (-96.3%)	0.46 (-42.4%)	0.80 (-43.5%)
Adam	0.34 (-80.6%)	0.04 (-99.5%)	0.29 (-63.9%)	0.29 (-79.5%)
AdamW	$5.33 \ (\gg 100\%)$	0.20 (-97.2%)	0.54 (-32.4%)	0.78 (-45.2%)
Signum	0.78 (-55.3%)	0.13 (-98.1%)	0.50 (-36.8%)	0.58 (-59.0%)
SignumW	$3185.69 \ (\gg 100\%)$	0.30 (-95.7%)	1.15 (+44.2%)	1.40 (-1.2%)

3.2 WEIGHT DECAY IS ESSENTIAL AND MOMENTUM ACCELERATES NC

Our experiments show that weight decay is necessary to reduce the NC metric across all optimizers and hyperparameter settings, as shown in Figure 3 for Signum and SGD, and earlier in Figure 1 for Adam and AdamW.

From the figures, we can conclude that larger weight decay leads to a stronger decrease of NC metrics. In particular, we show that adaptive optimizers with decoupled weight decay have much larger NC metrics, which are strictly away from zero, showing no sign of NC. In addition, we show empirically that momentum amplifies the effect of weight decay on the decrease of NC metrics in SGD, as shown in the heatmap in Figure 4. This implies that one achieves a decrease in the NC metrics both by increasing weight decay for fixed momentum or by increasing momentum for fixed non-zero weight decay. The effect of momentum on the NC metrics becomes larger for larger values of weight decay.

The experimental results are complemented by Theorem 3.1 showing that NC0 converges to 0 with an exponential rate trained with SGD, which is proportional to momentum and weight decay,

¹see Appendix C.5 for an introduction to UFM.

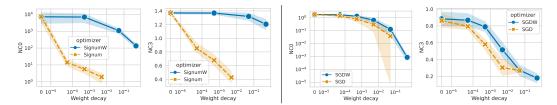


Figure 3: NC0 and NC3 metrics at the end of training for a ResNet9 trained on FashionMNIST for Signum and SignumW (left side) and SGD and SGDW (right side). Shaded area refers to one standard deviation across all trainings run with corresponding optimizer.

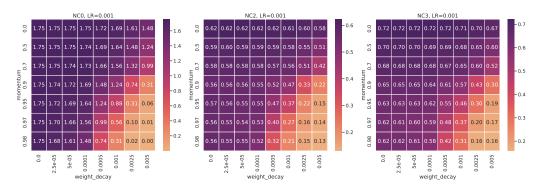


Figure 4: Heatmap of NC0, NC2 and NC3 for varying values of momentum and weight decay on ResNet9 trained on FashionMNIST with SGD.

highlighting that NC cannot be achieved without weight decay and that momentum accelerates the convergence of NC metrics.

Note that Theorem 3.1 holds for any model $f(\mathbf{W}, \theta, x) = \mathbf{W}h_{\theta}(x)$ with the last layer being a linear classifier and with any backbone h_{θ} parameterized by θ .

Theorem 3.1 (SGD promotes NC0). Assume a model of the form $f(\mathbf{W}, \theta, x) = \mathbf{W} h_{\theta}(x)$ is trained using cross-entropy loss with stochastic gradient descent (SGD) and momentum $\beta \in [0, 1)$, weight decay $\lambda \in [0, 1)$, and learning rate $\eta > 0$ sufficiently small. The last-layer weights \mathbf{W} are updated according to:

$$\begin{aligned} \mathbf{V}_{t+1} &= \beta \mathbf{V}_t + \nabla_{\mathbf{W}_t} L_{\text{CE}} + \lambda \mathbf{W}_t, \\ \mathbf{W}_{t+1} &= \mathbf{W}_t - \eta \mathbf{V}_{t+1}. \end{aligned}$$

Let $\alpha_t := \frac{1}{K} \|\mathbf{W}_t^{\top} \mathbf{1}\|_2^2$ denote the NCO metric. Suppose the condition $\frac{2\lambda}{\log \beta^{-1}} < 1$ holds. Fix a finite time horizon T > 0 such that $\eta t \le T$.

Then there exists an absolute constant C > 1 such that

$$\alpha_t \le \begin{cases} C \exp\left(-\frac{\lambda \eta t}{\log \beta^{-1}}\right) + \mathcal{O}(\eta) & \text{if } \beta > 0, \\ C(1 - \lambda \eta)^t & \text{if } \beta = 0, \end{cases}$$

Proof. The key observation is that the row sum of the loss gradient $\nabla L_{\text{CE}}(\mathbf{W}_t)^{\top} \mathbf{1}_K$ is zero, which largely simplifies the NC0 metric to only be dependent on the weight decay λ and momentum β . For the details of the proof, please refer to Subsection E in the Appendix.

Remark 3.2 (Convergence rate). Note that existing theoretical analysis on NC usually focuses on unconstrained feature models (UFM), and on the analysis of global solutions and the optimization landscape. This theorem goes beyond landscape analysis and computes the actual convergence rate, which allows us to quantify the effect of momentum and WD on NC.

3.3 WEIGHT DECAY COUPLING MATTERS

While weight decay has been theoretically shown to be essential for NC in prior works (Pan & Cao, 2024; Jacot et al., 2022), these works ignore how weight decay is applied by treating L_2 -regularization of the gradient and applying weight decay directly on parameters as equivalent. However, we note that this equivalency only holds for vanilla SGD and not for adaptive optimizers, such as Adam or AdamW, nor when momentum is applied. In particular, our experiments reveal that NC does not emerge under SignumW and AdamW under realistic settings. This highlights the crucial role of coupled weight decay – that is L_2 -regularization applied directly within the gradient update – as a requirement for NC. This subtle yet important distinction has been largely overlooked in prior literature.

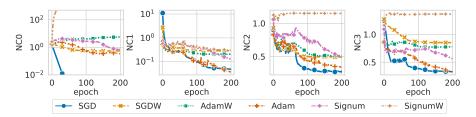


Figure 5: NC metrics throughout training on a ResNet9 trained on FashionMNIST.

Importantly, tracking the evolution of the NC metrics (Figure 5) and the singular values of centered class means M and the last-layer weight W (Figure 6) throughout training (here shown for a ResNet9 trained on FashionMNIST), one can see that using adaptive optimizers with decoupled weight decay leads to fundamentally different dynamics of the NC metrics and singular values despite all models reaching TPT, where training error is (almost) zero.

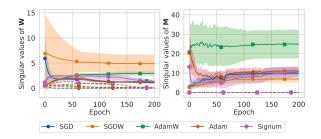


Figure 6: Singular values of last-layer weights \mathbf{W} (left) and centered class means \mathbf{M} (right) throughout training. The dotted line corresponds to the smallest singular value and the full line corresponds to the average singular value, excluding the smallest singular value. Singular values for SignumW are out-of-range and are shown in Figure 16 in the appendix.

Specifically, Figure 6 shows that the smallest singular value of **W** increases during training with AdamW and SignumW, indicating failure to satisfy NC3. Additionally, NC0 and the nonzero singular values of **M** grow throughout training and exhibit high variance, suggesting that NC2 is also less well-fulfilled in these settings.

In Figure 5, we further observe that SGD and Adam achieve the lowest NC metric values, while AdamW, SignumW, and SGDW saturate early at much higher levels. Although the NC metrics for Signum are slightly larger than for SGD and Adam, they continue to decrease over time, suggesting potential convergence to NC under longer training.

Finally, our experiments in Figure 1 and Figure 3 demonstrate that the NC0 and NC3 metrics of AdamW and SignumW remain significantly larger than those of Adam and Signum, even when using weight decay several orders of magnitudes higher. This indicates that models trained with AdamW or SignumW are consistently farther from achieving NC. Interestingly, the NC metrics for SGD and SGDW remain relatively close – only slightly shifted – showing that the gap between coupled and decoupled weight decay has a more pronounced effect in adaptive optimizers than in SGD. This suggests the effect is not simply due to greater weight decay accumulation through momentum but stems from a deeper interaction with the optimization dynamics.

3.4 INTERPOLATING ADAMW AND ADAM

To further investigate why AdamW fails to exhibit neural collapse (NC) while Adam does, we conducted an ablation study by "interpolating" between the two optimizers. Specifically, we implemented a variant that combines both coupled weight decay (as in Adam) and decoupled weight decay (as in AdamW). For each run, we varied the strength of the coupled weight decay while adjusting the decoupled component such that the total weight decay remained fixed at 0.0005. The momentum was set to 0.9 across all configurations.

As shown in Figure 7, increasing the coupled component leads to a smooth improvement in NC metrics—particularly NC0, NC2, and NC3—while the validation accuracy remains largely unaffected. This experiment suggests that coupled weight decay is a critical factor in enabling neural collapse, yet it is not strictly necessary for achieving strong generalization performance, as all configurations yield similar validation accuracy. This strengthens a point raised earlier about the limitations of NC to understand generalization Hui et al. (2022).

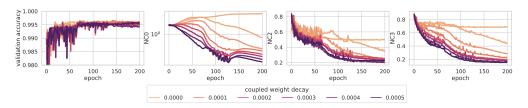


Figure 7: Interpolating Adam and AdamW by varying the coupled and decoupled weight decay. Total weight decay was fixed to 0.0005. Note that coupled weight decay = 0 is equivalent to AdamW and coupled weight decay = 0.0005 is equivalent to Adam. Experiments trained on ResNet9 with MNIST.

This observation is supported by our theoretical results in Theorem 3.3 and Theorem 3.4, which show that SignGD with decoupled weight decay fails to satisfy NC0 and therefore cannot converge to a neural collapse solution, whereas SignGD with coupled weight decay exhibits different behaviour. We note that SignGD corresponds to a special case of Adam and AdamW when the parameters β_1 , β_2 , and ε are set to zero.

Theorem 3.3 (Sign GD with decoupled weight decay avoids NC0). Consider sign GD with (decoupled) weight decay $\lambda > 0$ and step size $\eta > 0$ on the UFM loss $L_{CE}(\mathbf{WH}, \mathbf{I}) = \sum_{n=1}^{N} L_{CE}(\mathbf{Wh}_n, \mathbf{e}_n)$, where the feature $\mathbf{H} = \mathbf{M}^*$ is fixed to an NC solution and only the weight \mathbf{W} is trained:

$$\mathbf{W}_{t+1} = \mathbf{W}_t - \eta(sign(\nabla_{\mathbf{W}_t} L_{CE}) + \lambda \mathbf{W}_t)$$

Define the NC0 metric $\alpha = \|\mathbf{W}_t^{\top} \mathbf{1}_K\|_2^2$ as before. Then we have

$$\lim_{t \to \infty} \alpha_t = \frac{(K-2)^2}{\lambda^2}.$$

In particular, α_t does not vanish as $t \to \infty$.

Proof idea: The key observation is that the signed loss gradient $sign(\nabla L_{CE}(\mathbf{W}_t))$ in this setting is constant in t, simplifying the following computation. See Appendix E for the full proof.

Theorem 3.4 (Sign GD with coupled weight decay can lead to NCO). Consider sign GD with (coupled) weight decay $\lambda > 0$ and step size $\eta > 0$ on the UFM loss $L_{CE}(\mathbf{WH}, \mathbf{I}) = \sum_{n=1}^{N} L_{CE}(\mathbf{Wh}_n, \mathbf{e}_n)$, where the feature $\mathbf{H} = \mathbf{M}^*$ is fixed to an NC solution and only the weight \mathbf{W} is trained:

$$\mathbf{W}_{t+1} = \mathbf{W}_t - \eta(sign(\nabla_{\mathbf{W}_t} L_{CE} + \lambda \mathbf{W}_t))$$

We initialize $\mathbf{W}_0 = 0 \in \mathbb{R}^{K \times K}$ and define the covariance matrix $\mathbf{C}_t = \mathbf{W}_t \mathbf{W}_t^{\top}$ and the scalar $\alpha_t = \langle \mathbf{C}_t, \hat{\mathbf{J}} \rangle_F$ where $\hat{\mathbf{J}} = \frac{1}{K} \mathbf{1} \mathbf{1}^{\top}$. Then there exists some $T_2, T_1 > 1$ such that α increases for $t \in [1, T_1]$, decreases for $t \in [T_1 + 1, T_2]$ and oscillates with range $\mathcal{O}\left(d^3\eta^2\right)$ for $t > T_2 + 1$.

Proof. See Appendix E.

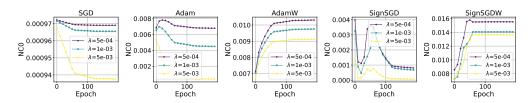


Figure 8: Training dynamic of NC0 with optimizers SGD, Adam, AdamW, AdamU ($\beta_1 = \beta_2 = 0$), AdamW0 ($\beta_1 = \beta_2 = 0$).

The key difference between the results of Theorem 3.3 and Theorem 3.4 lies in how coupled weight decay affects the signed gradient during training. As the weight norm $\|\mathbf{W}\|$ increases, the coupled decay term can eventually flip the sign of the gradient, altering the trajectory of the NC0 metric α_t . Initially, α_t grows at a similar rate in both cases, but their behaviors diverge once the decay term becomes dominant.

To illustrate this effect, we conducted a small-scale experiment using a simple MLP on a separable dataset with various optimizers. As shown in Figure 8, SignSGD displays non-monotonic dynamics in α_t , while SignSGDW exhibits steady convergence to a positive value. Similar patterns appear in Adam and AdamW, though more smoothed due to their adaptive updates.

4 DISCUSSION AND LIMITATIONS

In this section, we discuss new insights, additional considerations and limitations from the main results in Section 3. Additionally, we explore potential follow-up research directions that could provide theoretical explanations or extend our experiments to broader settings.

4.1 Interpreting NC Metrics in Practice

While NC is defined by the convergence of all NC metrics to zero in the limit, practical experiments never achieve exact zeros. Since NC is inherently a continuous rather than discrete phenomenon, it becomes necessary to define what constitutes the presence of NC in practice. This important issue has not been thoroughly addressed in the existing literature.

A further complication is that different NC metrics operate on different scales and these scales vary across settings of architectures and datasets. For example, in our experiments, the smallest observed values for NC2 and NC3 are on the order of 0.1, whereas NC1 can reach values an order of magnitude smaller.

In this work, we therefore refer to the emergence of NC in terms of relative strength. Specifically, we use the NC metric values at initialization as a baseline for models that do not exhibit NC, and use the smallest values achieved across all experiments as a reference point for models that do. This framing allows us to discuss the strength of NC emergence across different optimizers and settings.

4.2 The Redundant NC4 property

Readers may notice that we omit NC4 from the results in Section 3. This is because we observed that NC4 is consistently satisfied whenever the training accuracy approaches 100%, regardless of whether the other NC metrics (NC1–NC3) exhibit collapse. As shown in Figure 44, NC4 is largely uncorrelated with the other metrics. To maintain a clearer and more focused presentation, we therefore exclude NC4 from our main analysis.

4.3 PARTIAL NEURAL COLLAPSE

Another subtlety we observe is what we term *partial neural collapse*. As shown in Table 2, AdamW can achieve minimal values for NC1 and NC2 among all optimizers, even while NC0 diverges and NC3 is not satisfied. This indicates that NC properties may not always emerge jointly, contrary to the

Table 2: Final NC metrics for the run with the smallest absolute NC3 metric and > 99% training accuracy for each optimizer. Lower values (\downarrow) indicate stronger neural collapse. Values in parentheses represent percentages relative to the metric at initialization. Hyperparameters used for each optimizer can be found in Table 4.

Optimizer	$NC0_{\downarrow}$	NC1↓	NC2↓	NC3↓
SGD	1.53e-05 (< −99.5%)	0.02 (< -99.5%)	0.19 (-75.8%)	0.13 (-90.9%)
SGDW	1.54e-04 (< -99.5%)	0.01 (< -99.5%)	0.15 (-81.7%)	0.10 (-92.7%)
Adam	0.12 (< -93.2%)	0.04 (-99.5%)	0.23 (-71.6%)	0.17 (-88.2%)
AdamW	8.09 (≫100%)	0.01 (< -99.5%)	0.14 (-82.1%)	0.49 (-65.1%)

original claim in Papyan et al. (2020). Understanding the theoretical conditions under which only a subset of NC properties holds remains an intriguing open question.

4.4 LIMITATIONS OF THEORETICAL SUPPORT

Our experiments on Adam and AdamW are conducted on realistic models and datasets, whereas our theoretical results (Theorem 3.3, Theorem 3.4) focus on a simplified setting: SignGD applied to the unconstrained feature model. While this restricted setup already demonstrates that AdamW fails to achieve NC, it does not fully capture the complexity of deep neural networks or adaptive optimizers in practice. Nevertheless, we believe our proof techniques could be extended to explain why Adam may lead to NC in more general settings. Moreover, our theoretical analysis is limited to the training dynamics of NC0, chosen for its analytical tractability and strong empirical correlation with other NC metrics. A full theoretical understanding of NC1–NC3 under realistic optimization dynamics remains an open challenge, and we leave this direction for future work.

4.5 FUTURE RESEARCH

Other than the topic we have discussed in the previous subsections, our findings also open other intriguing avenues for future research.

- Empirical studies should be expanded to include larger models, such as Vision Transformers (ViTs) and DenseNets, as well as more diverse datasets, to assess the broader generality of our findings. Our preliminary results on ViT are available in Appendix D.4.5, and largely confirm our findings also extend to Transformers.
- Due to computational constraints, our study only analyzed NC properties in the last layer.
 However, previous works (Masarczyk et al., 2023; Rangamani et al., 2023) suggest that these properties may also manifest in intermediate layers. Investigating NC behavior across different depths could provide further insights into hierarchical feature representations.
- In addition to the optimizers (SGD, Adam, AdamW, Signum) studied in this work, novel first-order methods such as Lion (Chen et al., 2023) and Mars (Yuan et al., 2024), and second-order methods, such as Shampoo (Gupta et al., 2018), SOAP (Vyas et al., 2024) and Muon (Jordan et al.) demonstrated promising improvements in convergence and generalization. However, their effects on NC remain largely unexplored.

5 CONCLUSION

In this paper we have conducted an extensive number of experiments to elucidate the role of the optimization algorithm in the emergence of the neural collapse (NC) phenomenon. In particular, our experiments consistently show that coupled weight decay is necessary for achieving small NC metrics. While the role of weight decay in the context of NC has been studied in the literature before, this is the first paper distinguishing between coupled and decoupled weight decay. Moreover, our theoretical results show that the resulting training dynamics differ considerably and one needs to take this into account. These findings underscore the limitations of existing theoretical frameworks, which have studied NC mainly under gradient flow or gradient descent, and highlight the need for further investigation into the interplay between optimizers and NC.

REFERENCES

- Mouïn Ben Ammar, Nacim Belkhir, Sebastian Popescu, Antoine Manzanera, and Gianni Franchi. Neco: Neural collapse based out-of-distribution detection, 2024.
- Tina Behnia, Ganesh Ramachandra Kini, Vala Vakilian, and Christos Thrampoulidis. On the implicit geometry of cross-entropy parameterizations for label-imbalanced data. In Francisco Ruiz, Jennifer Dy, and Jan-Willem van de Meent (eds.), *Proceedings of The 26th International Conference on Artificial Intelligence and Statistics*, volume 206 of *Proceedings of Machine Learning Research*, pp. 10815–10838. PMLR, 25–27 Apr 2023.
- Jeremy Bernstein, Yu-Xiang Wang, Kamyar Azizzadenesheli, and Animashree Anandkumar. signsgd: Compressed optimisation for non-convex problems. In *International Conference on Machine Learning*, pp. 560–569. PMLR, 2018.
- Xiangning Chen, Chen Liang, Da Huang, Esteban Real, Kaiyuan Wang, Hieu Pham, Xuanyi Dong, Thang Luong, Cho-Jui Hsieh, Yifeng Lu, and Quoc V Le. Symbolic discovery of optimization algorithms. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine (eds.), *Advances in Neural Information Processing Systems*, volume 36, pp. 49205–49233. Curran Associates, Inc., 2023.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pp. 248–255. Ieee, 2009.
- Tomer Galanti, András György, and Marcus Hutter. On the role of neural collapse in transfer learning. *arXiv preprint arXiv:2112.15121*, 2021.
- Connall Garrod and Jonathan P. Keating. The persistence of neural collapse despite low-rank bias: An analytic perspective through unconstrained features, 2024.
- Vineet Gupta, Tomer Koren, and Yoram Singer. Shampoo: Preconditioned stochastic tensor optimization. In Jennifer Dy and Andreas Krause (eds.), Proceedings of the 35th International Conference on Machine Learning, volume 80 of Proceedings of Machine Learning Research, pp. 1842–1850. PMLR, 10–15 Jul 2018.
- Ernst Hairer, Gerhard Wanner, and Syvert P. Nørsett. *Solving Ordinary Differential Equations I Nonstiff Problems*. Springer-Verlag, Berlin,, 1993.
- X. Y. Han, Vardan Papyan, and David L. Donoho. Neural collapse under mse loss: Proximity to and dynamics on the central path, 2022.
- Md Yousuf Harun, Jhair Gallardo, and Christopher Kanan. Controlling neural collapse enhances out-of-distribution detection and transfer learning, 2025.
- Like Hui, Mikhail Belkin, and Preetum Nakkiran. Limitations of neural collapse for understanding generalization in deep learning. *arXiv* preprint arXiv:2202.08384, 2022.
- Arthur Jacot, François Ged, Berfin Şimşek, Clément Hongler, and Franck Gabriel. Saddle-to-saddle dynamics in deep linear networks: Small initialization training, symmetry, and sparsity, 2022.
- Arthur Jacot, Peter Súkeník, Zihan Wang, and Marco Mondelli. Wide neural networks trained with weight decay provably exhibit neural collapse, 2024.
- Jiachen Jiang, Jinxin Zhou, Peng Wang, Qing Qu, Dustin Mixon, Chong You, and Zhihui Zhu. Generalized neural collapse for a large number of classes. *arXiv preprint arXiv:2310.05351*, 2023.
- K Jordan, Y Jin, V Boza, Y Jiacheng, F Cecista, L Newhouse, and J Bernstein. Muon: An optimizer for hidden layers in neural networks, 2024b. *URL https://kellerjordan.github.io/posts/muon*.
- Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint* arXiv:1412.6980, 2014.
- Vignesh Kothapalli. Neural collapse: A review on modelling principles and generalization, 2023.

- Litian Liu and Yao Qin. Detecting out-of-distribution through the lens of neural collapse. *arXiv* preprint arXiv:2311.01479, 2023.
- Xuantong Liu, Jianfeng Zhang, Tianyang Hu, He Cao, Yuan Yao, and Lujia Pan. Inducing neural collapse in deep long-tailed learning. In *International conference on artificial intelligence and statistics*, pp. 11534–11544. PMLR, 2023.
 - Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization, 2019.
- Wojciech Masarczyk, Mateusz Ostaszewski, Ehsan Imani, Razvan Pascanu, Piotr Miłoś, and Tomasz Trzciński. The tunnel effect: Building data representations in deep neural networks. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine (eds.), *Advances in Neural Information Processing Systems*, volume 36, pp. 76772–76805. Curran Associates, Inc., 2023.
- Dustin G Mixon, Hans Parshall, and Jianzong Pi. Neural collapse with unconstrained features. *Sampling Theory, Signal Processing, and Data Analysis*, 20(2):11, 2022.
- Kaouther Mouheb, Marawan Elbatel, Stefan Klein, and Esther E Bron. Evaluating the fairness of neural collapse in medical image classification. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pp. 286–296. Springer, 2024.
- Leyan Pan and Xinyuan Cao. Towards understanding neural collapse: The effects of batch normalization and weight decay, 2024.
- Vardan Papyan, XY Han, and David L Donoho. Prevalence of neural collapse during the terminal phase of deep learning training. *Proceedings of the National Academy of Sciences*, 117(40): 24652–24663, 2020.
- Akshay Rangamani, Marius Lindegaard, Tomer Galanti, and Tomaso A Poggio. Feature learning in deep classifiers through intermediate neural collapse. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.), *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pp. 28729–28745. PMLR, 23–29 Jul 2023.
- Christos Thrampoulidis, Ganesh Ramachandra Kini, Vala Vakilian, and Tina Behnia. Imbalance trouble: Revisiting neural-collapse geometry. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neural Information Processing Systems*, volume 35, pp. 27225–27238. Curran Associates, Inc., 2022.
- Nikhil Vyas, Depen Morwani, Rosie Zhao, Mujin Kwun, Itai Shapira, David Brandfonbrener, Lucas Janson, and Sham Kakade. Soap: Improving and stabilizing shampoo using adam. *arXiv* preprint *arXiv*:2409.11321, 2024.
- Robert Wu and Vardan Papyan. Linguistic collapse: Neural collapse in (large) language models. *arXiv preprint arXiv:2405.17767*, 2024.
- Huizhuo Yuan, Yifeng Liu, Shuang Wu, Xun Zhou, and Quanquan Gu. Mars: Unleashing the power of variance reduction for training large models, 2024.
- Zhihui Zhu, Tianyu Ding, Jinxin Zhou, Xiao Li, Chong You, Jeremias Sulam, and Qing Qu. A geometric analysis of neural collapse with unconstrained features. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan (eds.), *Advances in Neural Information Processing Systems*, volume 34, pp. 29820–29834. Curran Associates, Inc., 2021.

Appendix

A LLM USAGE STATEMENT

We disclaim that we have used Large Language Models to refine a few sentences and additionally as a proxy of a search engine to retrieve additional related work.

The appendix is organized as follows. In Section B, we formally define the neural collapse (NC) phenomenon and introduce the metrics used in the experiments presented in the main text. In Appendix C, we review prior works related to our paper. Section D provides detailed descriptions and additional observations from our experiments. In Section E, we present the full proof of the theorems stated in the main text.

B NC METRICS

Neural collapse (NC), discovered by Papyan et al. (2020), is a striking phenomenon observed during the terminal phase of training (TPT) deep neural networks (DNN) for multi-class classification tasks, particularly when trained with cross-entropy (CE) loss. Formally, let the (trained) last-layer features of the DNN be denoted by \mathbf{h}_n , and concatenate them into a matrix $\mathbf{H} \in \mathbb{R}^{p \times N}$, where p is the width of the last layer and N is the number of training samples indexed by n. The output logits of the network are then computed as $\mathbf{W}_L\mathbf{H} + \mathbf{b}\mathbf{1}_N^\top \in \mathbb{R}^{K \times N}$, where $\mathbf{W}_L \in \mathbb{R}^{K \times p}$ is the last-layer weight, $\mathbf{b} \in \mathbb{R}^K$ is the bias vector, and K is the number of classes.

The DNN is trained using the CE loss computed on the logits:

$$CE(\mathbf{W}_{L}, \mathbf{H}) = -\sum_{n=1}^{N} \log \left(\frac{\exp(\mathbf{W}_{L} \mathbf{h}_{n})_{y_{n}}}{\sum_{k=1}^{K} \exp(\mathbf{W}_{L} \mathbf{h}_{n})_{k}} \right),$$

where $y_n \in [K]$ denotes the class label index of the feature vector \mathbf{h}_n . Let $\mathcal{C}_k \stackrel{\text{def.}}{=} n \in [N] : y_n = k$ be the index set of data points belonging to class $k \in [K]$. In this paper, we assume that the classes are balanced, i.e., $|\mathcal{C}_k|$ is equal for all $k \in [K]$. For the effects of class imbalance on NC, we refer the reader to Han et al. (2022); Thrampoulidis et al. (2022); Behnia et al. (2023).

Let $\mu_k \stackrel{\text{def.}}{=} \frac{1}{|\mathcal{C}_k|} \sum_{n \in \mathcal{C}_k} \mathbf{h}_n$ be the class mean for each class k. The global mean of all classes is given by $\mu_G = \frac{1}{K} \sum_{k=1}^K \mu_k$ and centered class means are defined as $\bar{\mu}_k = \mu_k - \mu_G$. Let the between-class covariance $\Sigma_B \in \mathbb{R}^{p \times p}$ and the within-class covariance $\Sigma_W \in \mathbb{R}^{p \times p}$ be:

$$\boldsymbol{\Sigma}_{B} = \frac{1}{K} \sum_{k=1}^{K} \bar{\boldsymbol{\mu}}_{k} \bar{\boldsymbol{\mu}}_{k}^{\top},$$

$$\boldsymbol{\Sigma}_{W} = \frac{1}{N} \sum_{k=1}^{K} (\mathbf{h}_{n} - \boldsymbol{\mu}_{G}) (\mathbf{h}_{n} - \boldsymbol{\mu}_{G})^{\top}.$$

We also concatenate the centered class means into a matrix $\mathbf{M} \stackrel{\text{def.}}{=} (\bar{\mu}_1,...,\bar{\mu}_K) \in \mathbb{R}^{p \times K}$.

With these definitions in place, we now conceptually outline the NC properties and introduce corresponding metrics to quantitatively measure these properties in our experiments.

NC1 - Variability Collapse The first property of neural collapse (NC1) describes the collapse of features to their respective class means. Formally, this means that the distance between a feature

²For simplicity, we interchangeably refer to an input $\mathbf{x} \in \mathbb{R}^d$ and its corresponding last-layer feature $\mathbf{h} \in \mathbb{R}^p$ after the parameters of the network have converged during TPT and the mapping $\mathbf{x} \mapsto \mathbf{h}$ is fixed.

vector \mathbf{h}_n and its corresponding class mean $\boldsymbol{\mu}_k$ approaches zero:

$$\|\mathbf{h}_n - \boldsymbol{\mu}_k\|_2 \to 0, \forall k \in [K], \ n \in \mathcal{C}_k.$$

A corresponding metric is defined as Zhu et al. (2021); Kothapalli (2023); Ammar et al. (2024):

$$\mathcal{NC}1 \stackrel{\text{def.}}{=} \frac{1}{K} \text{Tr}[\mathbf{\Sigma}_W \mathbf{\Sigma}_B^{\dagger}] \tag{2}$$

where † denotes the Moore-Penrose pseudo-inverse.

NC2 - Convergence of Class Means to Simplex ETF The second property of neural collapse (NC2) describes the convergence of class means to a simplex equiangular tight frame (ETF), where the angles between the means are maximally symmetric. Formally, this property can be expressed as:

$$\begin{cases} \|\bar{\boldsymbol{\mu}}_j\|_2 - \|\bar{\boldsymbol{\mu}}_k\|_2 & \to 0 \\ \left\langle \frac{\bar{\boldsymbol{\mu}}_j}{\|\bar{\boldsymbol{\mu}}_j\|_2}, \frac{\bar{\boldsymbol{\mu}}_k}{\|\bar{\boldsymbol{\mu}}_k\|_2} \right\rangle & \to \frac{K}{K-1} \delta_{jk} - \frac{1}{K-1}, \forall j, k \in [K]. \end{cases}$$

To measure this property, we define two metrics capturing the equinormality and equiangularity of the centered class means Papyan et al. (2020); Ammar et al. (2024):

$$\mathcal{NC}2_n = \frac{\text{std}_k\{\|\bar{\mu}_k\|_2\}}{\text{avg}_k\{\|\bar{\mu}_k\|_2\}};$$
(3)

$$\mathcal{NC}2_a = \operatorname{avg}_{k \neq k'} \left| \left\langle \frac{\bar{\boldsymbol{\mu}}_k}{\|\bar{\boldsymbol{\mu}}_k\|_2}, \frac{\bar{\boldsymbol{\mu}}_{k'}}{\|\bar{\boldsymbol{\mu}}_{k'}\|_2} \right\rangle + \frac{1}{K - 1} \right|. \tag{4}$$

Here, $std_{\bullet}(\cdot)$ and $avg_{\bullet}(\cdot)$ denote the standard deviation and mean, respectively, over the specified index.

An alternative metric for NC2, introduced by Kothapalli (2023), directly measures the deviation of the centered class means from a simplex ETF:

$$\mathcal{NC}2 \stackrel{\text{def.}}{=} \frac{1}{K^2} \left\| \frac{\mathbf{M}^{\top} \mathbf{M}}{\|\mathbf{M}^{\top} \mathbf{M}\|_F} - \mathbf{M}^* \right\|_F$$
 (5)

where

$$\mathbf{M}^* \stackrel{\text{def.}}{=} \frac{1}{\sqrt{K-1}} \left(\mathbf{I}_K - \frac{1}{K} \mathbf{J}_K \right),$$

 $\mathbf{I}_K \in \mathbb{R}^{K \times K}$ is the identity matrix and $\mathbf{J} \in \mathbb{R}^{K \times K}$ is the matrix of ones. Note that $\mathcal{NC}2_n, \mathcal{NC}2_a \to 0 \iff \mathcal{NC}2 \to 0$.

NC2W - Convergence of Weight Rows to Simplex ETF In addition to NC2, we define a related property, NC2W, which describes the convergence of the rows of the last-layer weights $\mathbf{W}_L \in \mathbb{R}^{K \times p}$ to a simplex ETF. If the third NC property, NC3 (described later), holds, then NC2 and NC2W are equivalent. However, to study partial NC, it is essential to decouple these properties and measure NC2 and NC2W separately.

To measure NC2W, Zhu et al. (2021) introduced the following metric:

$$\mathcal{NC}2\mathcal{W} \stackrel{\text{def.}}{=} \frac{1}{K^2} \left\| \frac{\mathbf{W}_L \mathbf{W}_L^{\top}}{\left\| \mathbf{W}_L \mathbf{W}_L^{\top} \right\|_F} - \mathbf{M}^* \right\|_F.$$
 (6)

While this metric measures the overall alignment of \mathbf{W}_L with a simplex ETF, it does not account for the equinormality and equiangularity of the rows of \mathbf{W}_L . To address this, we introduce the following metrics:

$$\mathcal{NC}2W_n = \frac{\operatorname{std}_k\{\|\mathbf{w}_k\|_2\}}{\operatorname{avg}_k\{\|\mathbf{w}_k\|_2\}}$$
(7)

$$\mathcal{NC}2W_a = \operatorname{avg}_{k \neq k'} \left| \left\langle \frac{\mathbf{w}_k}{\|\mathbf{w}_k\|_2}, \frac{\mathbf{w}_{k'}}{\|\mathbf{w}_{k'}\|_2} \right\rangle + \frac{1}{K - 1} \right|$$
(8)

where $\mathbf{w}_k^{\top} \in \mathbb{R}^p$ is the k-th row of \mathbf{W}_L .

NC2M - Convergence of Product to Simplex ETF Finally, Zhu et al. (2021); Kothapalli (2023) proposed a metric that interpolates between NC2 and NC2W: ³

$$\mathcal{NC}2\mathcal{M} \stackrel{\text{def.}}{=} \frac{1}{K^2} \left\| \frac{\mathbf{W}_L \mathbf{M}}{\|\mathbf{W}_L \mathbf{M}\|_F} - \mathbf{M}^* \right\|_F.$$
 (9)

Note that $\mathcal{NC}2, \mathcal{NC}2\mathcal{W} \to 0 \implies \mathcal{NC}2\mathcal{M} \to 0$ but the converse does not hold.

NC3 - Convergence to Self-Duality The third property of neural collapse (NC3) describes that the rows of the last-layer weight align with the column of the class means, that is,

$$\left\| \frac{\mathbf{W}_L}{\|\mathbf{W}_L\|_F} - \frac{\mathbf{M}^\top}{\|\mathbf{M}^\top\|_F} \right\|_F \to 0;$$

the corresponding metric is an obvious one Papyan et al. (2020); Garrod & Keating (2024):

$$\mathcal{NC3} \stackrel{\text{def.}}{=} \frac{1}{Kp} \left\| \frac{\mathbf{W}_L}{\|\mathbf{W}_L\|_E} - \frac{\mathbf{M}^\top}{\|\mathbf{M}^\top\|_E} \right\|_E$$
 (10)

NC4 - Simplification of Nearest-Class-Center (NCC) The fourth property of neural collapse (NC4) describes that the classifier decision boundaries become equivalent to those derived by a nearest-class-mean classifier, that is,

$$\operatorname*{arg\,max}_{k}\langle\mathbf{w}_{k},\mathbf{h}\rangle\rightarrow\operatorname*{arg\,min}\left\Vert \mathbf{h}-\boldsymbol{\mu}_{k}\right\Vert _{2}$$

for any test feature $\mathbf{h} \in \mathbb{R}^p$; hence we can fix a test set of features $\{\mathbf{h}_n^{\text{test}}\}_{n=1}^{N^{\text{test}}}$ define the metric:

$$\mathcal{NC4} \stackrel{\text{def.}}{=} \frac{1}{N^{\text{test}}} \sum_{n=1}^{N^{\text{test}}} \mathbf{1} \{ \arg \max_{k} \langle \mathbf{w}_k, \mathbf{h}_n^{\text{test}} \rangle = \arg \min_{k} \|\mathbf{h}_n^{\text{test}} - \boldsymbol{\mu}_k\|_2 \}$$
(11)

where 1 is the indicator function.

The above NC properties hold if their corresponding metrics approach zero (except for NC4, which approach one) as the training step $t \to \infty$. A solution \mathbf{W}_L , \mathbf{H} satisfying these properties is referred to as an NC solution.

To observe the interpolation between partial and full NC, we introduce a weaker property:

NC0 - Zero Row Sum of Last-Layer Weight This new property describes that the rows of the last-layer weight W_L sums up to zero with the corresponding metric

$$\mathcal{NC}0 \stackrel{\text{def.}}{=} \frac{1}{n} \left\| \mathbf{W}_{L}^{\mathsf{T}} \mathbf{1} \right\|_{2}, \tag{12}$$

Note that $\mathcal{NC}2\mathcal{W} \to 0 \implies \mathcal{NC}0 \to 0$ but the converse does not hold.

The analogous property for the last-layer features, **Zero Column Sum of Last-Layer Features**, holds automatically because the columns of M are centered class means:

$$\sum_{k=1}^{K} \bar{\mu}_k = \sum_{k=1}^{K} (\mu_k - \mu_G) = 0.$$

Thus, NC0 for the last-layer weights already represents a form of duality similar to NC3.

³In the original works, this metric was used to evaluate self-duality. However, in this paper, we decouple the NC properties to study the effects of implicit biases on each individually.

C ADDITIONAL RELATED WORK

C.1 WEIGHT DECAY AND NEURAL COLLAPSE

Weight Decay has been shown to be essential for NC in prior works, like (Zhu et al., 2021; Pan & Cao, 2024; Jacot et al., 2024). However, their statements on weight decay are for (quasi-)optimal solutions in oversimplified models, which ignore the complex interaction between non-convex loss landscape and optimizers. Please see Appendix C.5 for an example.

C.2 EMPIRICAL STUDIES ON THE EMERGENCE OF NEURAL COLLAPSE

Neural collapse has also been studied beyond the original problem setting, which assumes few balanced classes as well as noise-free labels. Notably, Wu & Papyan (2024) studied the occurrence of NC for large language models, which do not satisfy any of the original assumption. Jiang et al. (2023) studied neural collapse for a large number of classes, while Mouheb et al. (2024) studied the influence of imbalanced in medical image classification on NC.

C.3 APPLICATIONS OF NEURAL COLLAPSE

The observation of neural collapse (NC) has inspired a growing body of follow-up work that applies NC metrics across various settings. In the context of out-of-distribution (OOD) detection, Ammar et al. (2024) propose a novel post-hoc detection method based on the geometric properties of NC, while Harun et al. (2025) show that explicitly controlling for NC1 can enhance OOD detection performance. Notably, the latter also claim that AdamW leads to NC, based on empirical results where NC3 values hover around 0.5 across different models—mirroring the misleading metrics reported in Table 2. As we demonstrate in the main text, however, this does not indicate true NC. This discrepancy underscores the need for a more precise and systematic framework for evaluating NC – one of the central contributions of this work.

In a separate line of inquiry, Liu et al. (2023) study the impact of class imbalance on NC and propose explicit feature regularization terms to induce NC under imbalanced distributions, resulting in improved model performance.

C.4 COUPLED WEIGHT DECAY IN THE CONTEXT OF NEURAL COLLAPSE

To the best of our knowledge, no prior work has investigated the role of optimizer choice in the context of NC. When minimizing the objective in Equation (1) or Equation (13), the weight decay induced by the L2-regularization parameter λ is coupled with the training loss. However, with the introduction of AdamW Loshchilov & Hutter (2019), decoupled weight decay has become the default in many modern optimizers. This paper aims to bridge this gap by systematically examining the impact of coupled versus decoupled weight decay on the emergence of NC.

C.5 Unconstrained Feature Model

The unconstrained feature model (UFM) Mixon et al. (2022); Zhu et al. (2021) is a simplified theoretical framework commonly used to study neural collapse (NC). In UFM, the last layer feature is replaced by a trainable matrix $\mathbf{H} = (\mathbf{h}_n)_{n=1}^N$, referred to as the *unconstrained feature*, which mimics the role of feature extraction layers in deep neural networks (DNN). For analytical simplicity, the layer following the unconstrained feature is often assumed to be linear \mathbf{W} , making UFM a special case of deep linear networks (DLN):

$$\min_{\mathbf{W}, \mathbf{H}} \sum_{n=1}^{N} \ell(\mathbf{W}\mathbf{h}_n, \mathbf{y}_n) + \frac{\lambda}{2} ||\mathbf{W}||^2 + \frac{\lambda}{2} ||\mathbf{H}||^2,$$
(13)

simplifying the minimization problem in Equation (1). In this paper, the loss ℓ is always assumed to be the cross-entropy (CE) loss, because it is the standard loss used in multi-classification tasks.

Zhu et al. (2021) has reported positive results on NC using UFM. Informally it holds that:

Theorem C.1 (Theorem 3.1 and 3.2 in Zhu et al. (2021)). Any global optimal solution of UFM is an NC solution, while all other critical points are strict saddles. As a result, for random initialization, it is almost surely that gradient descent finds an NC solution.

Zhu et al. (2021) also experimented NC on realistic models with optimizers like SGD and Adam, concluding the universality of NC across different optimizers.

D EXPERIMENT

The experiments of this work, particularly regarding computing the NC metrics, were based on code in Wu & Papyan (2024), which can be found at Github repository https://github.com/rhubarbwu/neural-collapse, which was published under the MIT license. The implementation of VGG9 was based on Code taken from https://github.com/jerett/PyTorch-CIFAR10. The author granted explicit permission to use the code.

An overview of the experiments that were conducted in this work can be found in Table 3, which resulted in a total number of 36 different experimental settings of (architecture × optimizer × dataset) combinations. Each optimizer optimizer was trained using three different learning rates, six different values of momentum and six different values of weight decay, resulting in 108 training runs per optimizer and 3.888 training runs in total. Some of the runs diverged or only achieved suboptimal training performance, which were then discarded. In total we had 2.500 "valid" training runs, which reached at least 99% training accuracy, which were considered for for the subsequent data analysis.

Table 3: Overview of experiments conducted in this work.

Architectures	Optimizers	Datasets		
ResNet9, VGG9	SGD, SGDW, Adam, AdamW, Signum, SignumW	MNIST, FashionMNIST, CIFAR10		

D.1 DETAILS ON CHOICE OF HYPERPARAMETERS

Every model was trained over 200 epochs with a batch size of 128. The learning rate λ was chosen to be in $\lambda \in \{0.001, 0.01, 0.0679\}$ for SGD and SGDW (the last learning rate was also reported in the original work by Papyan et al. (2020)) and $\lambda \in \{0.001, 0.005, 0.01\}$ for Adam, AdamW, Signum, and SignumW because most trainings diverged with larger learning rates during initial experimental training runs. The learning rate was decayed by a factor of 10 after one third and two third of training as has been done in original work by Papyan et al. (2020). Momentum μ (or β_1 for Adam, AdamW, Signum, and SignumW) was chosen to be in the range $\mu \in \{0, 0.5, 0.7, 0.9, 0.95, 0.98\}$ for all optimizers and weight decay WD was chosen to be in the range WD $\in \{0, 5e^{-5}, 5e^{-4}, 5e^{-3}, 0.05, 0.5\}$ for SGD, SGDW, Adam, and Signum and WD $\in \{0, 5e^{-4}, 0.05, 0.5, 5, 10\}$ for SignumW and AdamW. The main motivation for using AdamW and Signum W with much larger weight decay values was based on the hypothesis that the effect of weight decay is reduced due to decoupling. The β_2 parameter in Adam and AdamW was left to its default value of 0.999.

D.2 DETAILS ON COMPUTATIONAL RESOURCES

All experiments, including preliminary experiments as well as the final 3.888 experiments were run on 5 NVIDIA RTX4090 GPUs with 24 GB RAM. Since the models and the batch size was comparably small, actually only 3 GB GPU memory per training was required. Each training took between 8 and 16 minutes, leading to a total of 500-1000 GPU hours of training.

Table 4: Hyperparameters for each optimizer to achieve the smallest NC3 metric shown in Table 2.

		7.5	
Optimizer	Learning rate	Momentum/ β_1	Weight decay
SGD	0.01	0.9	0.05
SGDW	0.0679	0.5	0.05
Adam	0.005	0.98	0.05
AdamW	0.005	0.95	5
Signum	0.001	0.9	0.05
SignumW	0.001	0.98	10

Table 5: Summary of regression fit between NC3 and NC0

Experiment	n	\hat{eta}	$SE(\hat{\beta})$	t-value	p-value	95 % CI	R^2 / Adj R^2	F-statistic
LR=0.001	170	0.1903	0.008	24.262	0.000	[0.175, 0.206]	0.778 / 0.777	588.6
LR=0.005	74	0.2017	0.012	16.252	0.000	[0.177, 0.226]	0.786 / 0.783	264.1
LR=0.01	114	0.1439	0.007	19.892	0.000	[0.13, 0.158]	0.779 / 0.777	395.7
LR=0.0679	41	0.1771	0.012	14.367	0.000	[0.152, 0.202]	0.841 / 0.837	206.4
all	399	0.1582	0.005	32.760	0.000	[0.149, 0.168]	0.730 / 0.729	1073

D.3 DETAILS ON REGRESSION FIT BETWEEN NC3 AND NC0

In this subsection we provide additional details regarding the regression fit between NC3 and NC0. For the sake of completeness, we show the regression fit in Figure 9 again below. In addition, we have also computed a regression fit across all training runs, which converged, and all learning rates, shown in Figure 10. A summary of the regression fit can be found in Table 5, showing that more than 70% of the variation in NC3 can be explained by NC0.

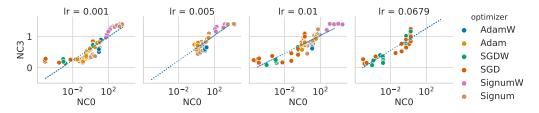


Figure 9: Figure 2 shown again for ease of reading. NC0 weakly correlates with NC3 across different optimizers and learning rates (here shown for ResNet9 trained on FashionMNIST).

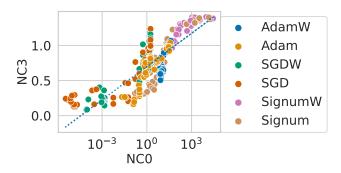


Figure 10: NC0 correlates with NC3 even when considered across all learning rates together (here shown for ResNet9 trained on FashionMNIST).

D.4 ADDITIONAL EXPERIMENTAL RESULTS

D.4.1 ABLATION STUDY ON TRAINING EPOCHS

As Neural collapse occurs at the terminal phase of training, it is natural to control for the effect that the number of training epochs has on the final NC metrics. After all, it is possible that the emergence of NC occurs at different speeds for different optimizers.

For this reason, we conducted two ablation studies, in which we prolong the training in two settings: We train a ResNet9 in FashionMNIST, which corresponds to the setting which is shown in Figure 1, for 2000 epochs with LR=0.0005 and momentum=0.9 for both optimizers. We note that in this setting, AdamW reaches 100% training accuracy already after around 700 epochs for all training runs with WD ≤ 0.05 . The results can be found in Figure 13 While this leads to some improvement of the final

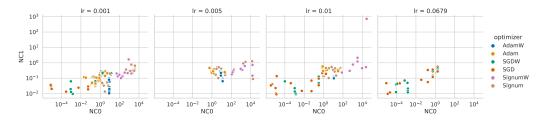


Figure 11: NC0 vs. NC1 across different optimizers and learning rates (here shown for ResNet9 trained on FashionMNIST).

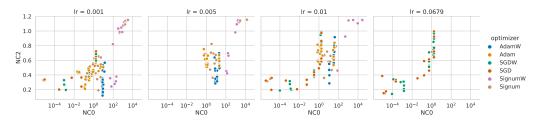


Figure 12: NC0 vs. NC2 across different optimizers and learning rates (here shown for ResNet9 trained on FashionMNIST).

NC1 and NC2 metric for AdamW for some values of weight decay, this has barely an effect on NC0 and NC3.

Furthermore we extend training to up to 2000 epochs for selected runs from Figure 4. Concretely, these runs trained with a LR of 0.001 and the following combination of WD and momentum (mom, WD) $\in \{(0,0), (0.97,5e^{-5}), (0,5e^{-4}), (0.9,5e^{-4}), (0.9,0), (0.95,0.0025)\}$, which corresponds to different parts in the heatmap. The results can be found in Figure 14. While one can observe a general decrease of the NC metrics in all cases, the overall trend for increasing weight decay remains unchanged. Both figures indicate that training the models considered in this work for 200 epochs is sufficient to draw the conclusions that we make about the necessity of coupled WD for the emergence of full NC.

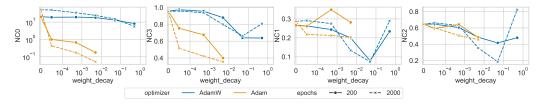


Figure 13: ResNet9 trained on FashionMNIST with Adam and AdamW for more epochs.

D.4.2 Unconstrained Feature Model

We also validated our results on the unconstrained feature model (UFM) (see Appendix C.5 for reference) with width d=512, K=10 classes and N=10.000 samples. The UFM was trained with Adam, AdamW and SGDMW with momentum=0.9 and varying $lr \in \{0.1, 0.3, 0.5, 1.0\}$ and weight decay ranging from 0.0 to 0.05. We then filtered the results, by only including models which achieved 100% training accuracy. The results in can be found in Figure 15. The plots show that the NC metrics, in particular NC0 and NC3 remain at least one magnitude of order larger than the same metrics for Adam and SGDMW, highlighting that AdamW converges to a different solution than Adam, which is not NC.

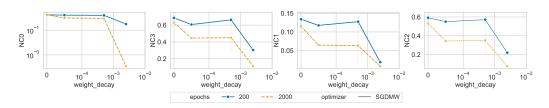


Figure 14: Selected runs from Figure 4 trained for more number of epochs.

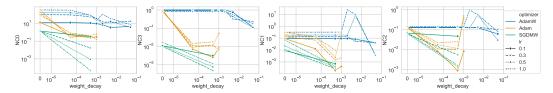


Figure 15: NC0 (left), NC3 (center left), NC1 (center right), and NC2 (right) for increasing weight decay.

D.4.3 MISSING PLOT: SINGULAR VALUE OF W AND M WITH SIGNUMW

The missing plot of the evolution of the singular values of the last-layer weights W and feature matrix M can be found in Figure 16.

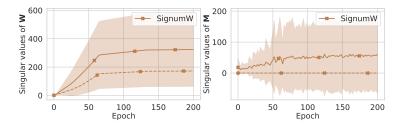


Figure 16: Singular values of last-layer weights W (left) and feature matrix M (right) throughout training for SignumW on ResNet9 trained on FashionMNIST. Dotted line corresponds do smallest singular value and full line corresponds to the average singular value excluding the smallest singular value.

D.4.4 COUPLED VS. DECOUPLED DECAY ON OTHER DATASETS

The comparison between coupled and decoupled decay on SGD, Adam, and Signum on other combinations of (architecture × dataset) can be found in the following pages below, which confirm our observations made earlier on the ResNet9 trained on FashionMNIST. While NC0 (visually) correlates well with NC3, it correlates considerably less with NC1 and NC2, although a general trend is still visible across all experiments.

ResNet50 on ImageNet1K We also conducted experiments on a ResNet50 trained on ImageNet1K Deng et al. (2009). The model was trained with Adam and AdamW for 90 epochs. We left out other optimizers due to limited resources. For both optimizers the learning rate was chosen as 0.0003 with a step-wise decay after 1/3 and 2/3 of training, momentum was chosen from $\{0.0, 0.5, 0.9\}$ and weight decay was chosen from $\{0.0, 1e^{-5}, 1e^{-4}, 1e^{-3}\}$. The resulting NC metrics can be found in Figure 17 and Figure 18, and confirm the conclusion that AdamW does not have full NC emergence.

VGG9 on FashionMNIST The comparison between coupled and decoupled weight decay on SGD, Adam, and Signum on a VGG9 trained on FashionMNIST can be found in Figure 19 and Figure 20. The relation between NC0 and NC3 can be found in Figure 23, between NC0 and NC1 in Figure 21, and between NC0 and NC2 in Figure 22.

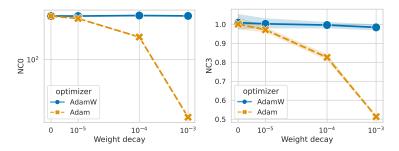


Figure 17: NC0 (left) and NC3 (right) metrics plotted against weight decay on a ResNet50 trained on ImageNet1K for Adam and AdamW. Shaded area refers to one standard deviation across all trainings run with corresponding optimizer.

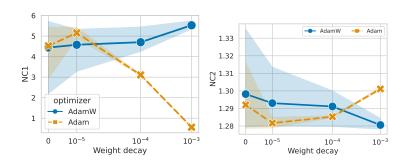


Figure 18: NC1 (left) and NC2 (right) metrics plotted against weight decay on a ResNet50 trained on ImageNet1K for Adam and AdamW. Shaded area refers to one standard deviation across all trainings run with corresponding optimizer.

ResNet9 on Cifar10 The comparison between coupled and decoupled weight decay on SGD, Adam, and Signum on a ResNet9 trained on Cifar10 can be found in Figure 24 and Figure 25. The relation between NC0 and NC3 can be found in Figure 28, between NC0 and NC1 in Figure 26, and between NC0 and NC2 in Figure 27.

VGG9 on Cifar10 The comparison between coupled and decoupled weight decay on SGD, Adam, and Signum can be found in Figure 29 and Figure 30. The relation between NC0 and NC3 can be found in Figure 33, between NC0 and NC1 in Figure 31, and between NC0 and NC2 in Figure 32.

ResNet9 on MNIST The comparison between coupled and decoupled weight decay on SGD, Adam, and Signum on a ResNet9 trained on MNIST can be found in Figure 34 and Figure 35. The relation between NC0 and NC3 can be found in Figure 38, between NC0 and NC1 in Figure 36, and between NC0 and NC2 in Figure 37.

VGG9 on MNIST The comparison between coupled and decoupled weight decay on SGD, Adam, and Signum on a VGG9 trained on MNIST can be found in Figure 39 and Figure 40. The relation between NC0 and NC3 can be found in Figure 43, between NC0 and NC1 in Figure 41, and between NC0 and NC2 in Figure 42.

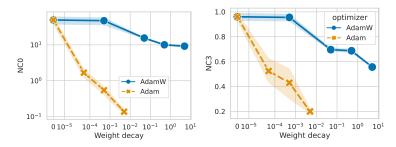


Figure 19: NC0 (left) and NC3 (right) metrics plotted against weight decay on a VGG9 trained on FashionMNIST for Adam and AdamW. Shaded area refers to one standard deviation across all trainings run with corresponding optimizer.

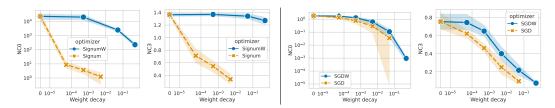


Figure 20: NC0 and NC3 metrics plotted against weight decay on a VGG9 trained on FashionMNIST for Signum and SignumW (left side) and SGD and SGDW (right side). Shaded area refers to one standard deviation across all trainings run with corresponding optimizer.

D.4.5 PRELIMINARY EXPERIMENTAL RESULTS ON VISION TRANSFORMER

We have also conducted preliminary experiments pretraining small Vision Transformers (ViT) on Cifar10 from scratch. Given that training ViTs is computationally much more expensive given the larger size of the model, we had to limit ourselves to a more restricted number of experiments. Specifically, we chose to train the ViT with Adam, AdamW, and SGD for 200 epochs with a batch size of 512 with momentum μ in the range $\mu \in \{0,0.8,0.9,0.95\}$ and weight decay WD $\in \{0,1e^{-5},1e^{-4},5e^{-4},1e^{-3},0.05,0.5\}$ for Adam and SGD and WD $\in \{0,1e^{-4},0.05,0.5,1,2,4\}$ for AdamW. We discarded all runs, which did not achieve a training accuracy of at least 50%. This mainly corresponded to training runs of SGD and Adam either with momentum=0 or WD ≥ 0.05 .

The ViT implementation is based on code from https://github.com/tintn/vision-transformer-from-scratch/tree/main, which is published under the MIT license. Specifically, the transformer model was chosen with a hidden dimension of 512, 6 hidden layers, and 8 attention heads, with no dropout applied.

Compared to the training procedure used in other settings, we employ a cosine-decay learning rate schedule with warm-up, where 5% of the total training steps are allocated to warm-up, and the base learning rate is set to 1×10^{-3} . Weight decay is applied to all layers except for LayerNorm and biases, which is standard practice.

The highest final test accuracy across all trainings was achieved by AdamW ($\beta_1 = 0.95$, WD = 0.5) with 83.67%, with a final test loss of 0.895. Notably, higher accuracy levels can be attained by increasing the network size and applying data augmentation or by using a pre-trained model as in Ammar et al. (2024). However, to ensure consistency with the experiments in the main study, we do not perform data augmentation due to limited computational resources. This likely explains the relatively lower test accuracy. Investigating the impact of data augmentation on the convergence to NC remains an interesting avenue for future work.

While we observe the general trend of decreasing NC metrics with increasing values of weight decay for SGD (Figure 45a), we note that in the case of ViTs the NC0 metric for both Adam and AdamW first increases before decreasing (Figure 45b, left), while the NC3 metric for both Adam and AdamW has a U-shape (Figure 45b, right). We also note that the ViT is much more sensitive to the choice of

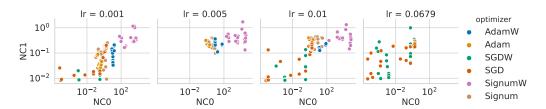


Figure 21: NC0 vs. NC1 on VGG9 trained on FashionMNIST. Note that the x-axis is plotted in log-scale.

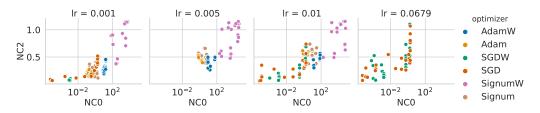


Figure 22: NC0 vs. NC2 on VGG9 trained on FashionMNIST. Note that the x-axis is plotted in log-scale.

weight decay and the training and validation accuracy degrades quickly due to overregularization, as can be seen in Figure 45c. A further investigation of these observations is left for future work.

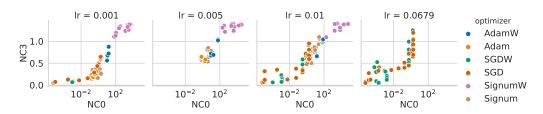


Figure 23: NC0 vs. NC3 on VGG9 trained on FashionMNIST. Note that the x-axis is plotted in log-scale.

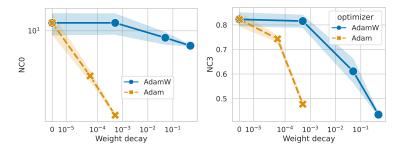


Figure 24: NC0 (left) and NC3 (right) metrics plotted against weight decay on a ResNet9 trained on Cifar10 for Adam and AdamW. Shaded area refers to one standard deviation across all trainings run with corresponding optimizer.

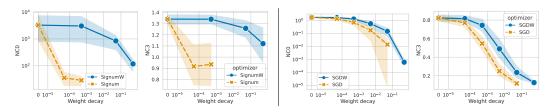


Figure 25: NC0 and NC3 metrics plotted against weight decay on a ResNet9 trained on Cifar10 for Signum and SignumW (left side) and SGD and SGDW (right side). Shaded area refers to one standard deviation across all trainings run with corresponding optimizer.

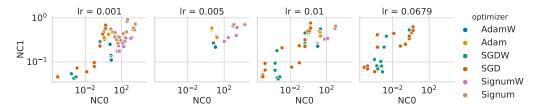


Figure 26: NC0 vs. NC1 on ResNet9 trained on Cifar10. Note that the x-axis is plotted in log-scale.

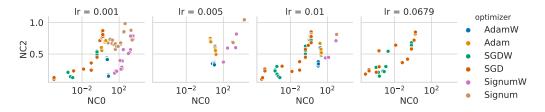


Figure 27: NC0 vs. NC2 on ResNet9 trained on Cifar10. Note that the x-axis is plotted in log-scale.

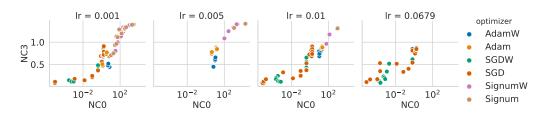


Figure 28: NC0 vs. NC3 on ResNet9 trained on Cifar10. Note that the x-axis is plotted in log-scale.

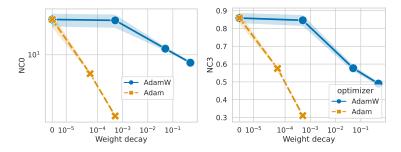


Figure 29: NC0 (left) and NC3 (right) metrics plotted against weight decay on a VGG9 trained on Cifar10 for Adam and AdamW. Shaded area refers to one standard deviation across all trainings run with corresponding optimizer.

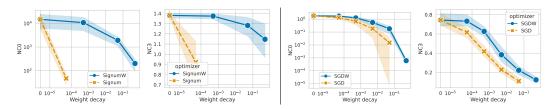


Figure 30: NC0 and NC3 metrics plotted against weight decay on a VGG9 trained on Cifar10 for Signum and SignumW (left side) and SGD and SGDW (right side). Shaded area refers to one standard deviation across all trainings run with corresponding optimizer.

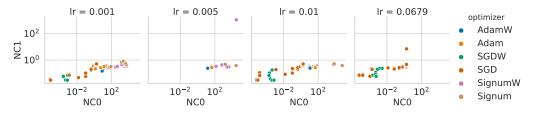


Figure 31: NC0 vs. NC1 on VGG9 trained on Cifar10. Note that the x-axis is plotted in log-scale.

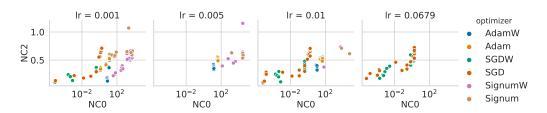


Figure 32: NC0 vs. NC2 on VGG9 trained on Cifar10. Note that the x-axis is plotted in log-scale.

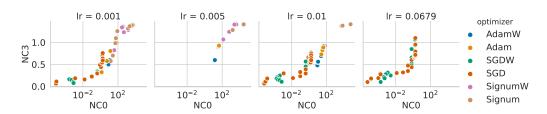


Figure 33: NC0 vs. NC3 on VGG9 trained on Cifar10. Note that the x-axis is plotted in log-scale.

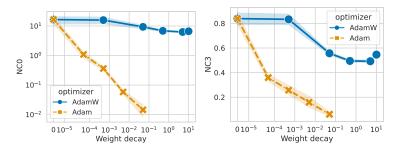


Figure 34: NC0 (left) and NC3 (right) metrics plotted against weight decay on a ResNet9 trained on MNIST for Adam and AdamW. Shaded area refers to one standard deviation across all trainings run with corresponding optimizer.

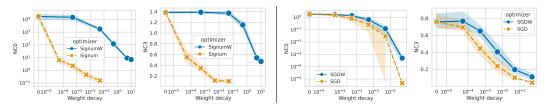


Figure 35: NC0 and NC3 metrics plotted against weight decay on a ResNet9 trained on MNIST for Signum and SignumW (left side) and SGD and SGDW (right side). Shaded area refers to one standard deviation across all trainings run with corresponding optimizer.

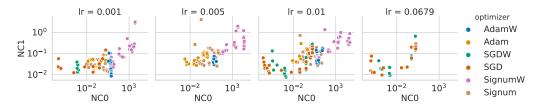


Figure 36: NC0 vs. NC1 on ResNet9 trained on MNIST. Note that the x-axis is plotted in log-scale.

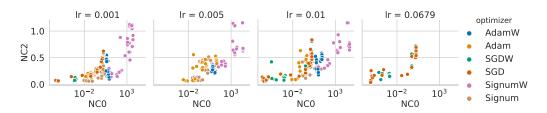


Figure 37: NC0 vs. NC2 on ResNet9 trained on MNIST. Note that the x-axis is plotted in log-scale.

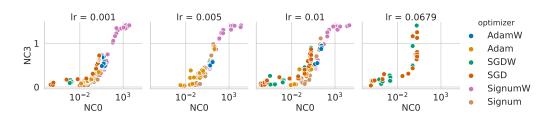


Figure 38: NC0 vs. NC3 on ResNet9 trained on MNIST. Note that the x-axis is plotted in log-scale.

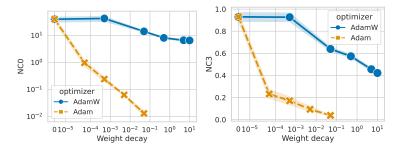


Figure 39: NC0 (left) and NC3 (right) metrics plotted against weight decay on a VGG9 trained on MNIST for Adam and AdamW. Shaded area refers to one standard deviation across all trainings run with corresponding optimizer.

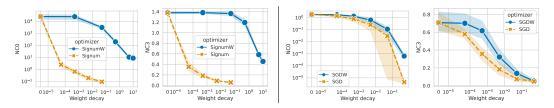


Figure 40: NC0 and NC3 metrics plotted against weight decay on a VGG9 trained on MNIST for Signum and SignumW (left side) and SGD and SGDW (right side). Shaded area refers to one standard deviation across all trainings run with corresponding optimizer.

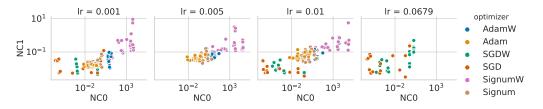


Figure 41: NC0 vs. NC1 on VGG9 trained on MNIST. Note that the x-axis is plotted in log-scale.

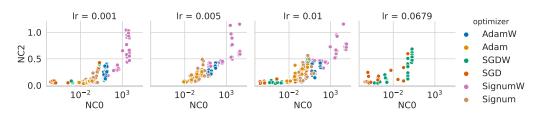


Figure 42: NC0 vs. NC2 on VGG9 trained on MNIST. Note that the x-axis is plotted in log-scale.

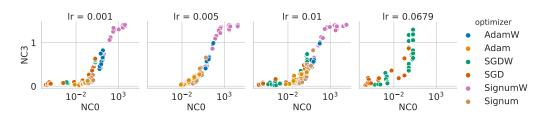


Figure 43: NC0 vs. NC3 on VGG9 trained on MNIST. Note that the x-axis is plotted in log-scale.

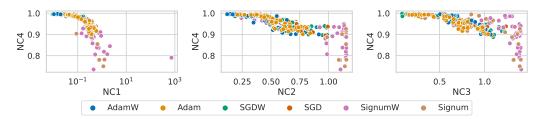
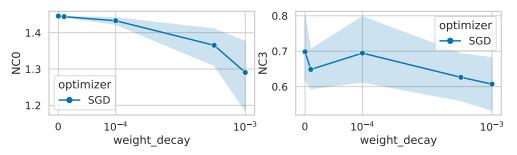
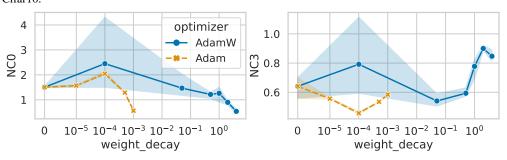


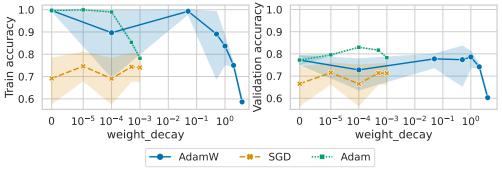
Figure 44: NC4 is largely uncorrelated with NC1-3 across different optimizers and learning rates.



(a) NC0 (left) and NC3 (right) metric for varying values of weight decay on a ViT trained with SGD on Cifar10.



(b) NC0 (left) and NC3 (right) metric for varying values of weight decay on a ViT trained with Adam and AdamW on Cifar10.



(c) Training accuracy (left) and validation accuracy (right) for varying values of weight decay on a ViT trained on Cifar10.

E Proofs

In this section, we will present the proof which is omitted in the main text.

Theorem E.1 (NC0 decay under momentum and weight decay). Assume a model of the form $f(\mathbf{W}, \theta, x) = \mathbf{W}h_{\theta}(x)$ is trained using cross-entropy loss with stochastic gradient descent (SGD) and momentum $\beta \in [0, 1)$, weight decay $\lambda \in [0, 1)$, and learning rate $\eta > 0$ sufficiently small. The last-layer weights \mathbf{W}_t are updated according to:

$$\mathbf{V}_{t+1} = \beta \mathbf{V}_t + \nabla_{\mathbf{W}_t} L_{\text{CE}} + \lambda \mathbf{W}_t,$$

$$\mathbf{W}_{t+1} = \mathbf{W}_t - \eta \mathbf{V}_{t+1}.$$
 (14)

Let $\alpha_t := \frac{1}{K} \|\mathbf{W}_t^{\top} \mathbf{1}\|_2^2$ denote the NCO metric. Suppose the condition $\frac{2\lambda}{\log \beta^{-1}} < 1$ holds. Fix a finite time horizon T > 0 such that $\eta t \le T$.

Then there exists an absolute constant C > 1 such that

$$\alpha_t \le \begin{cases} C \exp\left(-\frac{\lambda \eta t}{\log \beta^{-1}}\right) + \mathcal{O}(\eta) & \text{if } \beta > 0, \\ C(1 - \lambda \eta)^t & \text{if } \beta = 0, \end{cases}$$

Proof. Assume $\beta > 0$. We start from Lemma E.5: Then we have the expression:

$$\frac{1}{\eta}(\alpha_{t+1} - \alpha_t) = -2\beta\omega_t - 2\gamma_t - 2\lambda\alpha_t + \eta\nu_{t+1}$$

where $\omega_t \stackrel{\text{def.}}{=} \langle \mathbf{V}_t \mathbf{W}_t^{\top}, \hat{\mathbf{J}} \rangle$, $\gamma_t \stackrel{\text{def.}}{=} \langle \mathbf{G}_t \mathbf{W}_t^{\top}, \hat{\mathbf{J}} \rangle$, $\nu_t \stackrel{\text{def.}}{=} \langle \mathbf{V}_t \mathbf{V}_t^{\top}, \hat{\mathbf{J}} \rangle$.

We assume that $\mathbf{G}_t = \nabla L_{\text{CE}}(\mathbf{W}_t)$ and by Lemma E.4, $\langle \mathbf{G}_t \mathbf{A}^\top, \hat{\mathbf{J}} \rangle = \langle \mathbf{A} \mathbf{G}_t^\top, \hat{\mathbf{J}} \rangle = 0$ for any matrix \mathbf{A} . We derive a recursive formula for $\omega_t = \langle \mathbf{V}_t \mathbf{W}_t^\top, \hat{\mathbf{J}} \rangle$:

$$\omega_{t} = \langle (\beta \mathbf{V}_{t-1} + \mathbf{G}_{t-1} + \lambda \mathbf{W}_{t-1}) (\mathbf{W}_{t-1} - \eta \mathbf{V}_{t})^{\top}, \hat{\mathbf{J}} \rangle$$

$$= \beta \langle \mathbf{V}_{t-1} \mathbf{W}_{t-1}^{\top}, \hat{\mathbf{J}} \rangle + \lambda \langle \mathbf{W}_{t-1} \mathbf{W}_{t-1}^{\top}, \hat{\mathbf{J}} \rangle + \langle (-\beta \eta \mathbf{V}_{t-1} + \lambda \eta \mathbf{W}_{t-1}) \mathbf{V}_{t}^{\top}, \hat{\mathbf{J}} \rangle$$

$$= \beta \omega_{t-1} + \lambda \alpha_{t-1} + \langle (-\beta \eta \mathbf{V}_{t-1} + \lambda \eta \mathbf{W}_{t-1}) (\beta \mathbf{V}_{t-1} + \mathbf{G}_{t-1} + \lambda \mathbf{W}_{t-1})^{\top}, \hat{\mathbf{J}} \rangle$$

$$= \beta \omega_{t-1} + \lambda \alpha_{t-1} - \beta \eta (\beta \nu_{t-1} + \lambda \omega_{t-1}) + \lambda \eta (\beta \omega_{t-1} + \lambda \alpha_{t-1})$$

$$= \beta \omega_{t-1} + \lambda (1 + \lambda \eta) \alpha_{t-1} - \beta^{2} \eta \nu_{t-1}.$$
(15)

By extending Eq. (15) recursively for ω_{τ} for each $\tau = t-1, t-2, ..., 1$, we have

$$\omega_t = \lambda (1 + \lambda \eta) \sum_{\tau=0}^{t-1} \beta^{t-1-\tau} \alpha_\tau - \beta^2 \eta \sum_{\tau=1}^{t-1} \beta^{t-1-\tau} \nu_\tau.$$
 (16)

Next, we derive a recursive formula for $\nu_{t+1} = \langle \mathbf{V}_{t+1} \mathbf{V}_{t+1}^{\top}, \hat{\mathbf{J}} \rangle$:

$$\nu_{t+1} = \langle (\beta \mathbf{V}_t + \mathbf{G}_t + \lambda \mathbf{W}_t)(\beta \mathbf{V}_t + \mathbf{G}_t + \lambda \mathbf{W}_t)^\top, \hat{\mathbf{J}} \rangle = \beta^2 \nu_t + 2\beta \lambda \omega_t + \lambda^2 \alpha_t$$
 (17)

By Eq. (16) and (17) and the facts that $\alpha_t, \nu_t \geq 0, \beta, \lambda, \eta \in (0, 1)$, we bound Eq. (23) from above

$$\frac{1}{1553} \frac{1}{\eta} (\alpha_{t+1} - \alpha_t) = -2\beta(1 - \lambda \eta)\omega_t - (2\lambda - \lambda^2 \eta)\alpha_t + \beta^2 \eta \nu_t$$
1555
$$= -2\beta(1 - \lambda \eta) \left(\lambda(1 + \lambda \eta) \sum_{\tau=0}^{t-1} \beta^{t-1-\tau} \alpha_\tau - \beta^2 \eta \sum_{\tau=1}^{t-1} \beta^{t-1-\tau} \nu_\tau \right) - (2\lambda - \lambda^2 \eta)\alpha_t + \beta^2 \eta \nu_t$$
1558
$$\leq -2 \left(\lambda \sum_{\tau=0}^{t-1} \beta^{t-\tau} \alpha_\tau - \beta^2 \eta \sum_{\tau=1}^{t-1} \beta^{t-\tau} \nu_\tau \right) - \lambda \alpha_t + \beta^2 \eta \nu_t$$
1560
$$= -2\lambda \sum_{\tau=0}^{t-1} \beta^{t-\tau} \alpha_\tau + 2\beta^2 \eta \sum_{\tau=1}^{t-1} \beta^{t-\tau} \nu_\tau - \lambda \alpha_t + \beta^2 \eta \nu_t$$
1563
$$\leq -\lambda \sum_{\tau=0}^{t} \beta^{t-\tau} \alpha_\tau + 2\beta^2 \eta \sum_{\tau=0}^{t-1} \beta^{t-\tau} \nu_\tau.$$
(18)

By abuse of notation, denote by $\alpha(t)$ with $\alpha(0)=\alpha_0$ the gradient flow of α_t by setting $\dot{\alpha}(t)=\lim_{\eta\to 0}\frac{1}{\eta}(\alpha_{t+1}-\alpha_t)$, then we have

$$\dot{\alpha}(t) \le -\lambda \left(\int_0^t \beta^{t-\tau} \alpha(\tau) d\tau \right)$$

where the residual term $2\beta^2\eta\sum_{\tau=0}^t\beta^{t-\tau}\nu_{\tau}$ vanishes as $\eta\to 0$. By Grönwall's Lemma, solving an ODE with Laplacian using Lemma E.6, given that $\frac{2\lambda}{\log\beta^{-1}}<1$, we have

$$\alpha(t) \le C\alpha_0 \exp\left(-\frac{\lambda}{\log \beta^{-1}}t\right)$$
 (19)

for some absolute constant C > 1.

 Finally, by approximation of Euler discretization (see Theorem 7.5 in Hairer et al. (1993)), we have

$$|\alpha_t - \alpha(t)| = \mathcal{O}(\eta) \tag{20}$$

as long as $\eta t \leq T$. Replacing the leading constant coefficient by larger constant C if necessary, we obtain the claim. For the case where $\beta=0$, Equation (18) is reduced to: $\frac{1}{\eta}(\alpha_{t+1}-\alpha_t)\leq -\lambda\alpha_t$. Standard argument would lead to the claimed result.

Note that the above Theorem holds for any model $f(\mathbf{W}, \theta, x) = \mathbf{W}h_{\theta}(x)$ with last layer as linear classifier and with any backbone h_{θ} parameterized by θ .

However, the dynamics of Adam is more complicated, hence we further restrict the setting to SignGD, a special case of Adam, training a UFM.

Here, we assume a balanced dataset with only one element in each class $k \in [K]$. It is obvious to extend our result to multiple elements per class. Hence the total input N = K is equal to the number of classes and the UFM loss can be written as

$$L_{\text{CE}}(\mathbf{WH}, \mathbf{I}) = \sum_{n=1}^{N} L_{\text{CE}}(\mathbf{Wh}_n, \mathbf{e}_n),$$

where we can decouple the regularization $\frac{\lambda}{2} \|\mathbf{W}\|^2 + \frac{\lambda}{2} \|\mathbf{H}\|^2$ into weight decay.

By Zhu et al. (2021), we know that the UFM

$$\min_{\mathbf{W},\mathbf{H}} \sum_{n=1}^{N} L_{\text{CE}}(\mathbf{W}\mathbf{h}_{n},\mathbf{y}_{n}) + \frac{\lambda}{2} \|\mathbf{W}\|^{2} + \frac{\lambda}{2} \|\mathbf{H}\|^{2},$$

has unique global minimum W, H and no strict saddle points. In particular, $H = UM^*$ for some orthogonal matrix $U \in O(P)$. To further simplify the analysis, we assume that P = N = K with $H = M^*$. Then we have the followings:

Theorem E.2. Consider sign GD with (decoupled) weight decay $\lambda > 0$ and step size $\eta > 0$ on the UFM loss

$$L_{\mathit{CE}}(\mathbf{WH}, \mathbf{I}) = \sum_{n=1}^{N} L_{\mathit{CE}}(\mathbf{Wh}_n, \mathbf{e}_n),$$

where the feature $\mathbf{H} = \mathbf{M}^*$ is fixed to an NC solution and only the weight \mathbf{W} is trained:

$$\mathbf{W}_{t+1} = \mathbf{W}_t - \eta(sign(\nabla_{\mathbf{W}_t} L_{CE}) + \lambda \mathbf{W}_t)$$

with initialization $\mathbf{W}_0 = 0 \in \mathbb{R}^{K \times K}$. We define the covariance matrix $\mathbf{C}_t = \mathbf{W}_t \mathbf{W}_t^{\top}$ and the scalar $\alpha_t = \langle \mathbf{C}_t, \hat{\mathbf{J}} \rangle_F$ where $\hat{\mathbf{J}} = \frac{1}{K} \mathbf{1} \mathbf{1}^{\top}$. Then we have

$$\lim_{t \to \infty} \alpha_t = \frac{(K-2)^2}{\lambda^2}.$$

In particular, α_t does not vanish as $t \to \infty$.

1620 Proof. By Lemma E.4, we have $\nabla L_{\text{CE}}(\mathbf{W}) = \frac{1}{N}(\mathbf{S} - \mathbf{Y})\mathbf{H}^{\top} = \frac{1}{N}(\mathbf{S} - \mathbf{I}) \cdot \frac{1}{\sqrt{K-1}}(\mathbf{I} - \frac{1}{K}\mathbf{J}) = \frac{1}{N\sqrt{K-1}}(\operatorname{softmax}(\mathbf{W}\mathbf{H}) - \mathbf{I})$ since $(\operatorname{softmax}(\mathbf{W}\mathbf{H}) - \mathbf{I})\mathbf{J} = 0$. Since softmax has range between 0 and 1, we have

$$sign\left(\nabla L_{CE}(\mathbf{WH})\right) = \mathbf{J} - 2\mathbf{I},$$

that is, the signed gradient is -1 on the diagonal and +1 elsewhere. Note that this holds for all $\mathbf{W} \in \mathbb{R}^{K \times K}$. The sign GD updates can hence be written as:

$$\mathbf{W}_{t+1} = \mathbf{W}_t - \eta \Big[\underbrace{\mathbf{J} - 2\mathbf{I}}_{\operatorname{sign}(\nabla_{\mathbf{W}_t} L_{\text{CE}})} + \lambda \mathbf{W}_t \Big].$$
(21)

Since sign $(\nabla L_{CE}(\mathbf{W}_t))$ is constant, the dynamics collapse onto a scalar w_t :

$$\mathbf{W}_t = w_t (\mathbf{J} - 2\mathbf{I}),$$

which has the following recursive form:

$$w_{t+1} = (1 - \eta \lambda)w_t - \eta, \quad w_0 = 0.$$

Solve it and obtain

$$w_t = -\frac{1}{\lambda} \left[1 - (1 - \eta \lambda)^t \right].$$

Recall the definition:

$$\mathbf{C}_t = \mathbf{W}_t \mathbf{W}_t^{\top} \quad \hat{\mathbf{J}} = \frac{1}{K} \mathbf{1} \mathbf{1}^{\top} \text{and} \quad \alpha_t = \langle \mathbf{C}_t, \hat{\mathbf{J}} \rangle_F.$$

Since $\|(\mathbf{J} - 2\mathbf{I})^{\top}\mathbf{1}\|^2 = (K - 2)^2 K$ and the factor of 1/K gives $(K - 2)^2$, we have

$$\alpha_t = (K-2)^2 w_t^2$$

Therefore

$$\alpha_t = (K - 2)^2 \left[-\frac{1}{\lambda} \left(1 - (1 - \eta \lambda)^t \right) \right]^2 = \frac{(K - 2)^2}{\lambda^2} \left[1 - \left(1 - \eta \lambda \right)^t \right]^2.$$

As
$$t \to \infty$$
, $(1 - \eta \lambda)^t \to 0$, so

$$\alpha_{\infty} = \frac{(K-2)^2}{\lambda^2}.$$

Theorem E.3. Consider sign GD with (coupled) weight decay $\lambda > 0$ and step size $\eta > 0$ on the UFM loss

$$L_{CE}(\mathbf{WH}, \mathbf{I}) = \sum_{n=1}^{N} L_{CE}(\mathbf{Wh}_n, \mathbf{e}_n),$$

where the feature $\mathbf{H} = \mathbf{M}^*$ is fixed to an NC solution and only the weight \mathbf{W} is trained:

$$\mathbf{W}_{t+1} = \mathbf{W}_t - \eta(sign(\nabla_{\mathbf{W}_t} L_{CE} + \lambda \mathbf{W}_t))$$

with initialization $\mathbf{W}_0 = 0 \in \mathbb{R}^{K \times K}$. We define the covariance matrix $\mathbf{C}_t = \mathbf{W}_t \mathbf{W}_t^{\top}$ and the scalar $\alpha_t = \langle \mathbf{C}_t, \hat{\mathbf{J}} \rangle_F$ where $\hat{\mathbf{J}} = \frac{1}{K} \mathbf{1} \mathbf{1}^{\top}$. Then there exists some $T_2, T_1 > 1$ such that α increases for $t \in [1, T_1]$, decreases for $t \in [T_1 + 1, T_2]$ and oscillates with range $\mathcal{O}\left(d^3\eta^2\right)$ for $t > T_2 + 1$.

Proof. Throughout the training, we apply mathematical induction on the structure of \mathbf{W}_t : for all t, there exists $a_t, b_t \geq 0$ such that

$$\mathbf{W}_t = (a_t + b_t)\mathbf{I} - b_t\mathbf{J}.$$

Note that for t=0, the signed gradient is the same as in the case with decoupled weight decay in Theorem 3.3:

$$\operatorname{sign}(\nabla_{\mathbf{W}_t} L_{\operatorname{CE}} + \lambda \mathbf{W}_t) = \operatorname{sign}(\nabla_{\mathbf{W}_0} L_{\operatorname{CE}}) = \operatorname{sign}(\operatorname{softmax}(0) - \mathbf{I}) = \mathbf{J} - 2\mathbf{I}.$$

Hence, $\mathbf{W}_1 = \eta(2\mathbf{I} - \mathbf{J})$ where $a_1 = b_1 = \eta$. Following the same argument, as long as the coupled weight decay $\lambda \mathbf{W}_t$ does not flip sign of the gradient, that is, for all $j \neq k$,

 $(\operatorname{softmax}(\mathbf{W}_t) - \mathbf{I})_{jk} + \lambda \mathbf{W}_{jk} > 0, \qquad (\operatorname{softmax}(\mathbf{W}_t) - \mathbf{I})_{jj} + \lambda \mathbf{W}_{jj} < 0, \tag{22}$

the sign gradient with coupled weight decay is still $\mathbf{J} - 2\mathbf{I}$, and thus

$$\mathbf{W}_{t+1} = (t+1)\eta(2\mathbf{I} - \mathbf{J}),$$

satisfying the induction hypothesis on \mathbf{W}_t with $a_t = b_t = \eta t$, as long as Equation (22) holds. Hence for all $t < T_1$, where $T_1 \in \mathbb{N}$ is chosen retroactively such that the conditions in Equation (22) hold, the conditions can be rewritten as:

$$\frac{e^{-b_t}}{e^{a_t} + (d-1)e^{-b_t}} - \lambda b_t > 0, \qquad \frac{-(d-1)e^{-b_t}}{e^{a_t} + (d-1)e^{-b_t}} + \lambda a_t < 0.$$

Write $\phi(t) = \frac{e^{-b_t}}{e^{a_t} + (d-1)e^{-b_t}} = \frac{1}{e^{(a_t+b_t)} + (d-1)}$, the conditions in Equation (22) can be written as:

$$\phi(t) > \lambda b_t$$
 $(d-1)\phi(t) > \lambda a_t$.

In particular, for all $t < T_1$, $a_t = b_t = \eta t$ and we have $\phi(t) = \frac{1}{e^{2\eta t} + (d-1)}$. Since $\phi(t) > 0$ decreases with t and $\lambda \eta t$ increases with t > 0 starting from zero, hence a minimum $T_1 \in \mathbb{N}$ must exist such that the first condition will break earlier than the second one:

$$\phi(T_1) = \frac{1}{e^{2\eta T_1} + (d-1)} < \lambda \eta T_1 < (d-1)\phi(T_1).$$

In which case, we have $sign(\nabla_{\mathbf{W}_{T_1}} L_{CE} + \lambda \mathbf{W}_{T_1}) = -\mathbf{J}$ and thus

$$\mathbf{W}_{T_1+1} = T_1 \eta(2\mathbf{I} - \mathbf{J}) + \eta \mathbf{J} = (a_t + b_t)\mathbf{I} - b_t \mathbf{J},$$

satisfying the induction hypothesis on W_{T_1+1} with

$$a_t = \eta t$$
, for $t \le T_1 + 1$; $b_t = \begin{cases} \eta t & \text{for } t \le T_1 \\ \eta(T_1 - 1), & \text{for } t = T_1 + 1 \end{cases}$

In particular, $a_{T_1+1} + b_{T_1+1} = 2\eta T_1$ so $\phi(T_1+1) = \phi(T_1)$. However, by the choice of T_1 ,

$$\phi(T_1+1) = \phi(T_1) > \lambda \eta(T_1-1)$$

and hence $(\nabla_{\mathbf{W}_{T_1+1}}L_{\mathrm{CE}} + \lambda \mathbf{W}_{T_1+1}) = \mathbf{J} - 2\mathbf{I}$, resulting in $a_{T_1+2} = \eta(T_1+2)$ and $b_{T_1+2} = \eta T_1$. Hence $\phi(T_1+2) < \phi(T_1+1) = \phi(T_1) < \lambda \eta T_1 = \lambda b_{T_1+2}$, breaking the first condition and we have $(\nabla_{\mathbf{W}_{T_1+2}}L_{\mathrm{CE}} + \lambda \mathbf{W}_{T_1+2}) = -\mathbf{J}$. Hence as long as the second condition holds, we have

$$a_t = \eta t, \ \, \text{for} \, t \leq T_1 + 1; \qquad b_t = \begin{cases} \eta t & \text{for} \, t \leq T_1 \\ \eta(T_1 - 1), & \text{for} \, t - T_1 > 0 \, \text{odd} \\ \eta T_1, & \text{for} \, t - T_1 > 0 \, \text{even} \end{cases}.$$

Hence, a_t will increase until the second condition breaks: there exists a minimum integer $T_2 > T_1$ such that

$$(d-1)\phi(T_2) = \frac{d-1}{e^{(a_t+b_t)} + d-1} < \frac{d-1}{e^{\eta(T_2+T_1-1)} + d-1} < \lambda a_{T_2} = \lambda \eta T_2.$$

Here we do two case distinctions: First assume the first condition also breaks, then we have $sign(\nabla_{\mathbf{W}_{T_2}} L_{CE} + \lambda \mathbf{W}_{T_2}) = 2\mathbf{I} - \mathbf{J}$ and thus

$$\mathbf{W}_{T_2+1} = \mathbf{W}_{T_2} - \eta(2\mathbf{I} - \mathbf{J}) = (a_t + b_t)\mathbf{I} - b_t\mathbf{J},$$

satisfying the induction hypothesis on W_{T_2+1} with

$$a_{T_2+1} = \eta(T_2-1) = a_{T_2-1};$$
 $b_{T_2+1} = \eta(T_1-1) = b_{T_2-1}$

Hence the later training will oscillate between $(a_t, b_t) = (\eta T_2, \eta T_1)$ and $(\eta (T_2 - 1), \eta (T_1 - 1))$, for all $t > T_2$, satisfying the induction hypothesis on \mathbf{W}_t .

Now we come to the other case where the first condition holds at $t=T_2$, then we have $sign(\nabla_{\mathbf{W}_{T_2}}L_{CE}+\lambda\mathbf{W}_{T_2})=\mathbf{J}$ and thus

 $\mathbf{W}_{T_2+1} = \mathbf{W}_{T_2} - \eta \mathbf{J} = (a_t + b_t)\mathbf{I} - b_t \mathbf{J},$

satisfying the induction hypothesis on \mathbf{W}_{T_2+1} with

$$a_{T_2+1} = \eta(T_2-1) = a_{T_2-1};$$
 $b_{T_2+1} = \eta(T_1-1+1) = b_{T_2}.$

Hence for $t = T_2 + 2$, the first condition breaks and the second condition holds, resulting in

$$a_{T_2+2} = \eta T_2 = a_{T_2};$$
 $b_{T_2+2} = \eta (T_1 - 1).$

Hence $\mathbf{W}_{T_2+3} = \mathbf{W}_{T_2+1}$ and thus the later training will oscillate between $(a_t, b_t) = (\eta(T_2-1), \eta T_1)$ and $(\eta T_2, \eta(T_1-1))$, for all $t > T_2$, satisfying the induction hypothesis on \mathbf{W}_t . Also, we have $a_t = \eta \min(t, T_2) \pm \eta$, $b_t = \eta \min(t, T_1) \pm \eta$ for all t. The remaining claims follow.

E.1 TECHNICAL LEMMATA

 Lemma E.4. Let $(\mathbf{X}, \mathbf{Y}) \in \mathbb{R}^{d \times N} \times \mathbb{R}^{K \times N}$ be a dataset where the labels \mathbf{Y} are written in columns of one-hot vectors. For each pair $(\mathbf{x}, \mathbf{y}) \in \mathbb{R}^D \times \mathbb{R}^K$, and a weight $\mathbf{W}_1 \in \mathbb{R}^{K \times d}$, define the cross-entropy as:

$$\ell(\mathbf{W}_1) \stackrel{\text{def.}}{=} -\sum_{k=1}^K \mathbf{y}_k \log \left(softmax(\mathbf{W}_1 \mathbf{x}) \right)_k = \log \left(1 + \sum_{k \neq y} \exp(\mathbf{w}_k - \mathbf{w}_y)^\top \mathbf{x}_i \right)$$

where $y = \arg \max_{k \in [K]} [\mathbf{y}]_k$ is the class index of \mathbf{x} . Let $\mathcal{L}_1(\mathbf{W}_1) = CE(\mathbf{W}_1\mathbf{X}, \mathbf{Y})$ be the average cross-entropy loss of the dataset (\mathbf{X}, \mathbf{Y}) . Then the loss gradient $\nabla \mathcal{L}_1(\mathbf{W}_1)$ is

$$abla \mathcal{L}_1(\mathbf{W}_1) = \frac{1}{N} (\mathbf{S} - \mathbf{Y}) \mathbf{X}^{ op}$$

where $\mathbf{S} = (\mathbf{s}_1, ... \mathbf{s}_N)$ and $\mathbf{s}_i = softmax(\mathbf{W}_1 \mathbf{x}_i)$ for each i. In particular, $\mathbf{1}_K^\top \nabla \mathcal{L}_1(\mathbf{W}_1) = 0$.

Proof. The expression of the loss gradient comes from simple calculus. The second statement comes from the fact that the L1 norms of a post-softmax vector and an one-hot vector are both equal to 1, that is,

$$\mathbf{1}_K^{\top} \mathbf{s}_i = \mathbf{1}_K^{\top} \mathbf{y}_i = 1 \, \forall i.$$

Lemma E.5. Assume the weight W_t is updated as follows:

$$\mathbf{V}_{t+1} = \beta \mathbf{V}_t + \mathbf{G}_t + \lambda \mathbf{W}_t$$
$$\mathbf{W}_{t+1} = \mathbf{W}_t - \eta \mathbf{V}_{t+1},$$

where G_t depends on W_t . Define

$$\alpha \stackrel{\text{def.}}{=} \frac{1}{K} \| \mathbf{W}_t^{\top} \mathbf{1} \|_2^2 \ge 0.$$

Then we have the expression:

$$\frac{1}{\eta}(\alpha_{t+1} - \alpha_t) = -2\beta\omega_t - 2\gamma_t - 2\lambda\alpha_t + \eta\nu_{t+1}$$

where $\omega_t \stackrel{\text{def.}}{=} \langle \mathbf{V}_t \mathbf{W}_t^{\top}, \hat{\mathbf{J}} \rangle$, $\gamma_t \stackrel{\text{def.}}{=} \langle \mathbf{G}_t \mathbf{W}_t^{\top}, \hat{\mathbf{J}} \rangle$, $\nu_t \stackrel{\text{def.}}{=} \langle \mathbf{V}_t \mathbf{V}_t^{\top}, \hat{\mathbf{J}} \rangle$.

Proof. Let $\mathbf{C}_t \stackrel{\text{def.}}{=} \mathbf{W}_t \mathbf{W}_t^{\top}$ be the covariance matrix. Notice that $\alpha_t = \langle \mathbf{C}_t, \hat{\mathbf{J}} \rangle$ where $\hat{\mathbf{J}} = \frac{1}{K} \mathbf{1} \mathbf{1}^{\top}$. By update rule of \mathbf{W}_t and \mathbf{V}_t :

$$\frac{1}{\eta} (\mathbf{C}_{t+1} - \mathbf{C}_t) = \frac{1}{\eta} \left((\mathbf{W}_t - \eta \mathbf{V}_{t+1}) (\mathbf{W}_t - \eta \mathbf{V}_{t+1})^\top - \mathbf{C}_t \right)$$
$$= -(\mathbf{V}_{t+1} \mathbf{W}_t^\top + \mathbf{W}_t \mathbf{V}_{t+1}^\top) + \eta \mathbf{V}_{t+1} \mathbf{V}_{t+1}^\top.$$

Applying the dot product $\langle \cdot, \hat{\mathbf{J}} \rangle_F$ on both sides, and denote $\omega_t \stackrel{\text{def.}}{=} \langle \mathbf{V}_t \mathbf{W}_t^\top, \hat{\mathbf{J}} \rangle$, $\gamma_t \stackrel{\text{def.}}{=} \langle \mathbf{G}_t \mathbf{W}_t^\top, \hat{\mathbf{J}} \rangle$, $\nu_t \stackrel{\text{def.}}{=} \langle \mathbf{V}_t \mathbf{V}_t^\top, \hat{\mathbf{J}} \rangle$, we have

$$\frac{1}{\eta}(\alpha_{t+1} - \alpha_t) = -2\langle \mathbf{V}_{t+1} \mathbf{W}_t^{\top}, \hat{\mathbf{J}} \rangle + \eta \langle \mathbf{V}_{t+1} \mathbf{V}_{t+1}^{\top}, \hat{\mathbf{J}} \rangle
= -2\langle (\beta \mathbf{V}_t + \mathbf{G}_t + \lambda \mathbf{W}_t) \mathbf{W}_t^{\top}, \hat{\mathbf{J}} \rangle + \eta \nu_{t+1}
= -2\beta \omega_t - 2\gamma_t - 2\lambda \alpha_t + \eta \nu_{t+1}$$
(23)

where in the first line we use the fact that $\hat{\bf J}$ is symmetric.

Lemma E.6. Assume $\lambda, \beta \in (0,1)$ such that $\frac{2\lambda}{\log \beta^{-1}} < 1$. The solution of the following ODE:

 $\dot{\alpha}(t) = -\lambda \left(\int_0^t \beta^{t-\tau} \alpha(\tau) d\tau \right) \tag{24}$

with initial condition $\alpha(0) = \alpha_0 > 0$ admits the following bound:

$$\alpha(t) \le C\alpha_0 \exp\left(-\frac{\lambda}{\log \beta^{-1}}t\right)$$

for some absolute constant C > 1.

Proof. Observe that we can write the integral in convolution:

$$\int_0^t \beta^{t-\tau} \alpha(\tau) d\tau = (\phi * \alpha)(t), \text{ where } \phi(t) = \beta^t.$$

Hence (24) can be written as

$$\dot{\alpha}(t) = -\lambda (\phi * \alpha)(t).$$

Let $\mathcal{L}\{\psi(t)\}(s)=\int_0^\infty e^{-st}\psi(t)dt$ denote the Laplace transform. Denote

$$\mathcal{A}(s) = \mathcal{L}\{\alpha(t)\}(s), \quad F(s) = \mathcal{L}\{\phi(t)\}(s).$$

Taking the Laplace transform of both sides:

$$\mathcal{L}\{\dot{\alpha}(t)\}(s) = -\lambda \mathcal{L}\{(\phi * \alpha)(t)\}(s). \tag{25}$$

And by integration by part and the property of convolution,

$$\mathcal{L}\{\dot{\alpha}(t)\}(s) = s\mathcal{A}(s) - \alpha(0)$$
 and $\mathcal{L}\{(\phi*\alpha)(t)\}(s) = F(s)\mathcal{A}(s)$.

Hence

$$s\mathcal{A}(s) - \alpha(0) = -\lambda F(s)\mathcal{A}(s).$$

Since $\beta^t = e^{(\log \beta)t}$, we get

$$F(s) = \mathcal{L}\{\beta^t\}(s) = \mathcal{L}\left\{e^{(\log \beta)t}\right\}(s) = \frac{1}{s - \log(\beta)} \quad \text{for } s > \log(\beta).$$

Substitute this back to Eq. (25) and we get:

$$s\mathcal{A}(s) - \alpha(0) = -\lambda \frac{1}{s - \log(\beta)} \mathcal{A}(s)$$

$$s\mathcal{A}(s) + \frac{\lambda}{s - \log(\beta)} \mathcal{A}(s) = \alpha(0)$$

$$\mathcal{A}(s) \left(\underbrace{s + \frac{\lambda}{s - \log(\beta)}}_{\frac{s^2 - s \log(\beta) + \lambda}{s - \log(\beta)}}\right) = \alpha(0)$$

$$\mathcal{A}(s) = \alpha(0) \cdot \underbrace{\frac{\left[s - \log(\beta)\right]}{s^2 - s\log(\beta) + \lambda}}_{(s-r_1)(s-r_2)}$$

1836
1837 where $r_1, r_2 = \frac{\log(\beta) \pm \sqrt{\left[\log(\beta)\right]^2 - 4\lambda}}{2}$. We do partial fractions and matching coefficients gives:

$$\frac{s - \log(\beta)}{(s - r_1)(s - r_2)} = \frac{A}{s - r_1} + \frac{B}{s - r_2} \implies A + B = 1, \quad -\log(\beta) = -Ar_2 - Br_1.$$

Since $r_1 + r_2 = \log(\beta)$, one finds

$$A = \frac{r_2}{r_2 - r_1}, \quad B = -\frac{r_1}{r_2 - r_1}.$$

Thus

$$\mathcal{A}(s) = \alpha(0) \left[\frac{r_2}{r_2 - r_1} \frac{1}{s - r_1} - \frac{r_1}{r_2 - r_1} \frac{1}{s - r_2} \right].$$

Recall the inverse of Laplacian transform: $\mathcal{L}^{-1}\left\{\frac{1}{s-r}\right\}(t) = e^{rt}$. Therefore,

$$\alpha(t) = \mathcal{L}^{-1}\{A(s)\}(t) = \alpha(0) \left[\frac{r_2}{r_2 - r_1} e^{r_1 t} - \frac{r_1}{r_2 - r_1} e^{r_2 t} \right].$$

Equivalently,

$$\alpha(t) = \alpha(0) \left[Ae^{r_1t} + Be^{r_2t} \right], \quad A = \frac{r_2}{r_2 - r_1}, B = -\frac{r_1}{r_2 - r_1},$$
 (26)

where

$$r_1, r_2 = \frac{\log(\beta) \pm \sqrt{\left[\log(\beta)\right]^2 - 4\lambda}}{2}$$

Since $\beta \in (0,1)$, set $L=-\log(\beta)>0$. By the first order approximation,

$$\sqrt{(\log \beta)^2 - 4\lambda} = \sqrt{L^2 - 4\lambda} = L - \frac{2\lambda}{L} + \mathcal{O}\left(\frac{\lambda^2}{L}\right)$$

Hence

$$r_1, r_2 = \frac{-L \pm \left(L - \frac{2\lambda}{L}\right)}{2} + \mathcal{O}\left(\frac{\lambda^2}{L}\right).$$

This gives:

$$r_1 = -\frac{\lambda}{L} + \mathcal{O}\left(\frac{\lambda^2}{L}\right), \quad r_2 = -L + \frac{\lambda}{L} + \mathcal{O}\left(\frac{\lambda^2}{L}\right).$$

Plugging r_1, r_2 into Eq. (26):

$$\alpha(t) \le C\alpha(0)e^{r_1t} = C\alpha(0)\exp\left(-\frac{\lambda}{L}t\right)$$

for some absolute constant C > 1. Plug in $L = -\log(\beta) = \log \beta^{-1}$ to finish the proof.