# Outliers with Opposing Signals Have an Outsized Effect on Neural Network Optimization

**Elan Rosenfeld**
Carnegie Mellon University
elan@cmu.edu

**Andrej Risteski**
Carnegie Mellon University

## Abstract

We identify a new phenomenon in network optimization which arises from the interaction of depth and a particular heavy-tailed structure in natural data. Our result offers intuitive explanations for several previously reported observations about network training dynamics. In particular, it implies a conceptually new cause for progressive sharpening and the edge of stability; we also highlight connections to other concepts in optimization and generalization including grokking, simplicity bias, and Sharpness-Aware Minimization.

Experimentally, we demonstrate the significant influence of paired groups of outliers in the training data with strong *opposing signals*: consistent, large magnitude features which dominate the network output throughout training and provide gradients which point in opposite directions. We describe how to identify these groups, explore what sets them apart, and carefully study their effect on the network's optimization and behavior. We complement these experiments with a mechanistic explanation on a toy example of opposing signals and a theoretical analysis of a two-layer linear network on a simple model. Our finding enables new qualitative predictions of training behavior which we confirm experimentally. It also provides a new lens through which to study and improve modern training practices for stochastic optimization, which we highlight via a case study of Adam versus SGD.

## 1 Introduction

There is a steadily growing list of intriguing properties of neural network (NN) optimization which are not readily explained by classical tools from optimization. Likewise, we have varying degrees of understanding of the mechanistic causes for each. Extensive efforts have led to possible explanations for the effectiveness of Adam [21], Batch Normalization [16] and other tools for successful training—but the evidence is not always entirely convincing, and there is certainly little theoretical understanding. These phenomena are typically considered in isolation—though they are not completely disparate, it is unknown what specific underlying causes they may share.

In this work, we identify a phenomenon in NN optimization which offers a new perspective on many of these prior observations and which we hope will contribute to a deeper understanding of how they may be connected. While we do not (and do not claim to) give a complete explanation, we present strong qualitative and quantitative evidence for a single high-level idea—one which naturally fits into several existing narratives and suggests a more coherent picture of their origin. Specifically, we demonstrate the prevalence of paired groups of outliers in natural data which have a significant influence on a network's optimization dynamics. These groups are characterized by the inclusion of one or more (relatively) large magnitude features that dominate the network's output at initialization and throughout most of training. In addition to their magnitude, the other distinctive property of these features is that they provide large, consistent, and *opposing* gradients, in that following one group's gradient to decrease its loss will increase the other's by a similar amount. Because of this structure, we refer to them
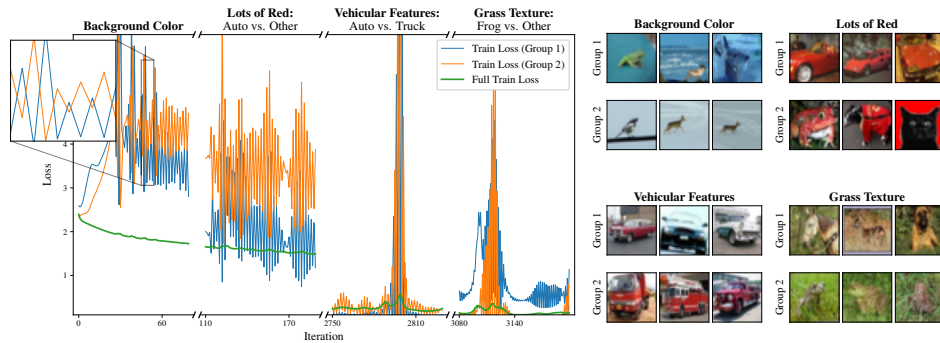
Figure 1: **Training dynamics of neural networks are heavily influenced by outliers with opposing signals.** We plot the overall loss of a ResNet-18 trained with GD on CIFAR-10, plus the losses of a small but representative set of outlier groups. These groups have consistent *opposing signals* (e.g., wheels and headlights can mean either `car` or `truck`). Throughout training, losses on these groups oscillate with growing and shrinking amplitude—this has an obvious correspondence to the intermittent spikes in overall loss and appears to be the direct cause of the edge of stability phenomenon.

as *Opposing Signals*. These features share a non-trivial correlation with the target task, but they are often not the "correct" (e.g., human-aligned) signal. In fact, in many cases these features perfectly encapsulate the classic statistical conundrum of "correlation vs. causation"—for example, a bright blue sky background does not determine the label of a CIFAR image, but it does most often occur in images of planes. Other features *are* very relevant, such as the presence of wheels and headlights in images of trucks and cars, or the fact that a colon often precedes either "the" or a newline token in written text.

Opposing signals are most easily understood with an example, which we will give along with a brief outline of their effect on training dynamics; a more detailed description is presented in Section 3. Fig. 1 depicts the training loss of a ResNet-18 trained with full-batch gradient descent (GD) on CIFAR-10, along with a few dominant outlier groups and their respective losses. In the early stages of training, the network enters a narrow valley in weight space which carefully balances the pairs' opposing gradients; subsequent sharpening of the loss landscape [19, 6] causes the network to oscillate with growing magnitude along particular axes, upsetting this balance. Returning to our example of a sky background, one step results in the class `plane` being assigned greater probability for all images with sky, and the next will reverse that effect. In essence, the "sky = `plane`" subnetwork grows and shrinks. The direct result of this oscillation is that the network's loss on images of planes with a sky background will alternate between sharply increasing and decreasing with growing amplitude, with the exact opposite occurring for images of *non*-planes with sky. As these pairs represent a small fraction of the data, this behavior is not immediately apparent from the overall training loss— but eventually, it progresses far enough that the overall loss spikes. As there is an obvious direct correspondence between these two events throughout, we conjecture that opposing signals are the direct cause of the *edge of stability* phenomenon [6].

We repeat this experiment across a range of vision architectures and training hyperparameters: though the precise groups and their order of appearance change, the pattern occurs consistently. We also verify this behavior for transformers on next-token prediction of natural text and small ReLU MLPs on simple 1D functions; we give some examples of opposing signals in text in Appendix C. However, we rely on images for exposition because it offers the clearest intuition. To isolate this effect, most of our experiments use GD, but we observe similar patterns during SGD which we present in Section 4.

**Summary of contributions.** The primary contribution of this paper is demonstrating the existence, pervasiveness, and large influence of opposing signals during NN optimization. We further present our current best understanding, with supporting experiments, of how these signals *cause* the observed training dynamics—in particular, we argue that it is a consequence of depth and steepest descent methods. We complement this discussion with a toy example and an analysis of a two-layer linear net on a simple model. Notably, though rudimentary, our explanation enables concrete qualitative predictions of NN behavior during training, which we confirm experimentally. It also provides a new lens through which to study modern stochastic optimization methods, which we highlight via a case study of SGD vs. Adam. We see possible connections between opposing signals and a wide variety of

phenomena in NN optimization and generalization, including *grokking* [42], *catapulting/slingshotting* [25, 50], *simplicity bias* [51], *double descent* [3, 34], and Sharpness-Aware Minimization [11]. We discuss these and other connections in Appendix A.

## 2   Characterizing and Identifying Opposing Signals

Though their influence on aggregate metrics is non-obvious, identifying outliers with opposing signals is straightforward. When training a network with GD, we track its loss on each individual training point. For a given iteration, we select the training points whose loss exhibited the most positive and most negative change in the preceding step (there is large overlap between these sets in successive steps). This set will sometimes contain multiple opposing signals, which we distinguish via visual inspection. This last detail means that the images we depict are not random, but we emphasize that it would not be correct to describe this process as cherry-picking: though precise quantification is difficult, these signals consistently obey the maxim "I know it when I see it". This is particularly true for images, such as the groups in Fig. 1 which have immediately recognizable patterns. To demonstrate this fact more generally, Appendix K contains the pre-inspection samples for a ResNet-18, VGG-11 [49], and a small Vision Transformer [9] at several training steps and for multiple seeds; we believe the implied groupings are immediate, even if not totally objective. We see algorithmic approaches to automatically clustering these samples as a direction for future study—for example, one could select samples by correlation in their loss time-series, or by gradient alignment.

**Measuring alternative metrics.**   Given how these samples are selected, several other characterizations seem appropriate. For instance, one-step loss change is often a reasonable proxy for gradient norm; we could also consider the largest eigenvalue of the loss of the *individual point*, or how much curvature it has in the direction of the overall loss's top eigenvector. For large networks these options are far more compute-intensive than our chosen method, but we can evaluate them on specific groups. In Fig. 29 in the Appendix we track these metrics for several opposing group pairs and find that they are consistently much larger than that of random samples from the training set.

**On the possibility of a formal definition.** Though the features and their exemplar samples are immediately recognizable, **we do not attempt to *exactly* define a "feature", nor an "outlier" with respect to that feature.** The presence of a particular feature is often ambiguous, and it is difficult to define a clear threshold for what makes a given point an outlier. Thus, instead of trying to exactly partition the data, we simply note that these heavy tails *exist* and we use the most obvious outliers as representatives for visualization. In Figs. 1 and 29 we choose an arbitrary cutoff of twenty samples per group. We also note that what qualifies as an opposing signal or outlier may vary over time. For visual clarity, Fig. 1 depicts the loss on only the most dominant group pair in its respective training phase, but this pattern occurs simultaneously for many different signals and at multiple scales throughout training. Further, the opposing signals are with respect to the model's internal representations (and the label), not the input space itself; this means that the definition is also a property of the architecture. In Fig. 30 in the Appendix we give a simple demonstration of this point, along with a more detailed discussion.

## 3   Understanding the Effect of Opposing Signals

Beyond noting their existence, our eventual goal will be to derive actionable insights from this finding. To do this, it is necessary to gain a better understanding of *how* these opposing signals lead to the observed behavior. Here we give a simplified "mental picture" which serves our current understanding this process: first a general discussion of why opposing signals are so influential, followed by a more mechanistic description with a toy example. This explanation is intentionally high-level, but we will eventually see how it gives concrete predictions of specific behaviors, which we then verify on real networks.

**Progressive sharpening, and why these features are so influential.**   In training a network to minimize predictive error, most information in the input will be unneeded—particularly with depth and high-dimensional inputs, only a small fraction will be propagated to the last linear layer [15]. Starting from random initialization, training a network aligns adjacent layers' singular values [47, 32] to amplify meaningful signal while downweighting noise,[1] growing *sensitivity* to the important

---

[1]In this discussion we use the term "noise" informally. We refer not necessarily to pure randomness, but more generally to input variation which is least useful in predicting the target.
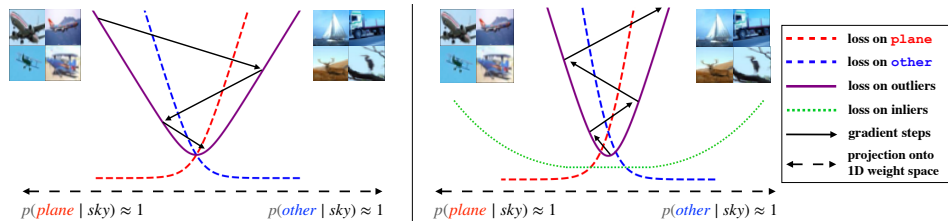
Figure 2: **A toy example illustrating the effect of opposing signals.** We project the loss to the hypothetical weight-space dimension "sky = plane". **Left:** Early optimization approaches the minimum, balancing the opposing gradients for *separate* losses plane and other. Progress continues through this valley, further growing the feature magnitude. **Right**: The valley sharpens and the iterates diverge. Because most images are less sensitive to this axis, the train loss is not noticeably affected at first. Eventually either (a) the outlier gradients' growth forces the network to downweight "sky", or (b) the iterates diverge enough that the weights "catapult" to a different basin.

signal. This sensitivity can be measured, for example, by the spectral norm of the input-output Jacobian, which grows during training [28]; it has also been connected to growth in the norm of the output layer [53]. Observe that as these norms grow, the network's sensitivity to changes in the *way* it processes inputs grows as well. Hypothetically, a small weight perturbation could massively increase loss by redirecting unhelpful noise to the subspace to which the network is most sensitive, or by changing how the last layer uses it. The increase of this sensitivity represents precisely the growth of loss Hessian spectrum, with the strength of this effect increasing with depth [52, 10, 32].[2]

Crucially, this sharpening also depends on the structure of the input. If the noise is independent of the target, it will be downweighted throughout training. In contrast, *genuine signals which oppose each other* will be retained and perhaps even further amplified by gradient descent; this is because the "correct" feature may be much smaller in magnitude (or not yet learned), so using the large, "incorrect" feature is often the most immediate way of minimizing loss. As a concrete example, observe that a randomly initialized network will lack the features required for the subtle task of distinguishing birds from planes. But it *will* capture the presence of sky, which is very useful for reducing loss on such images by predicting the conditional $p(\text{class} \mid \text{sky})$ (this is akin to the "linear-first" behavior described by Nakkiran et al. [33]). Thus, any method attempting to minimize loss as fast as possible (e.g., steepest descent) may actually upweight these features. Furthermore, amplified opposing signals will cause greater sharpening than random noise, because using a signal to the benefit of one group is maximally harmful for the other—e.g., confidently predicting plane whenever there is sky will cause enormous loss on images of other classes with sky. Since random noise is more diffuse, this effect is less pronounced. This description is somewhat abstract. To gain a more mechanistic understanding, we illustrate the precise dynamics on a toy example.

**Illustrating with a hypothetical example of gradient descent.** Consider the global loss landscape of a neural network: this is the function which describes how the loss changes as we move through parameter space. Suppose we identify a direction in this space which corresponds to the network's use of the "sky" feature to predict plane versus some other class. That is, we will imagine that whenever the input image includes a bright blue background, moving the parameters in one direction increases the logit of the plane class and decreases the others, and vice-versa. We will also decompose this loss—we consider separately the loss on images of planes and the loss on all other images. Fig. 2 depicts this heavily simplified scenario. Early in training, optimizing this network with GD will rapidly move towards the minimum along this direction. In particular, until the network learns to use more relevant signal, the direction of steepest descent will lead to a network which predicts the likelihood $p(\text{class} \mid \text{sky})$ whenever sky is present. Eventually, the gradient will no longer be dominated by this direction and will instead point "through the valley" [56]. However, until the network separates out the "sky" feature, this will simultaneously cause the sky feature to grow in magnitude. It will also cause an increase in the potential influence of this feature were the linear head to be selectively perturbed. Both these factors cause further sharpening. Continued optimization will oscillate across the minimum with growing magnitude, but this growth may not be immediately

---

[2]The coincident growth of these two measures was previously noted by Ma and Ying [28], Gamba et al. [13], MacDonald et al. [30], though they did not make explicit this connection to how the network processes different types of input variance.
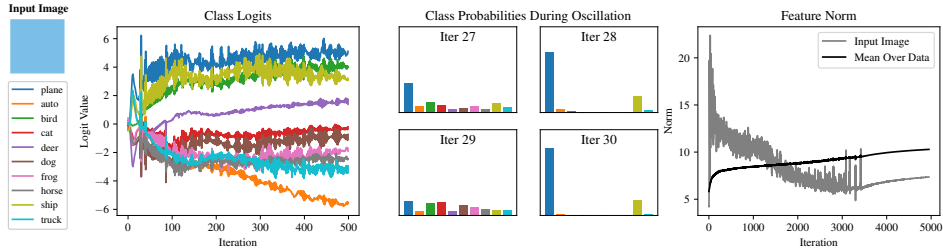
Figure 3: **Passing a sky-colored block through a ResNet during GD precisely tracks the predictions of our toy example. Left:** In the first phase, the network rapidly learns to use the sky, amplifying the feature and sharpening the loss. **Middle:** During oscillation, gradient steps alternate along the axis "sky = `plane`" (and a bit `ship`). **Right:** Oscillation originally amplifies the sky input; the network then slowly downweights this feature and learns to use other signal.

apparent. Furthermore, *progress orthogonal to these oscillations need not be affected*—we find some evidence that these two processes occur somewhat independently, which we present in Section 4. Returning to the loss decomposition, note that in addition to increasing in average magnitude, these oscillations will cause the losses to grow and *alternate*, with one group having high loss and then the other. Eventually the outliers' loss increases sufficiently and the overall loss spikes, either flattening the valley and returning to the first phase, or "catapulting" to a different basin [54, 25, 50]. This is the phenomenon depicted in Fig. 1.

Though this explanation lacks precise details, it does enable concrete predictions of network behavior during training. Fig. 3 tracks the predictions of a ResNet-18 on a synthetic image with all bright blue pixels. We see exactly the described behavior—initial convergence to the minimum along with rapid growth in feature norm, followed by oscillation in class probabilities. Over time, the network learns to use other signal and downweights the sky feature. We reproduce this figure for other inputs and for a VGG-11-BN [49] in the Appendix, with similar findings.

**Theoretical analysis of opposing signals in a simple model.** To demonstrate this effect more concretely, we study misspecified linear regression on inputs $x \in \mathbb{R}^d$ with a two-layer linear network. Though this model is quite simplified, it enables preliminary insight into the most important factors for these dynamics to occur. However, the concept of a "partially useful" signal seems to require a somewhat more complex model to properly capture (e.g., multinomial logistic regression), so we view this analysis only as an early investigation. Due to space constraints, we present our results in Appendix I. We study the trajectory of gradient flow, proving first an initial decrease in sharpness due to the model downweighting the noise, followed by a continuous, steady growth in sharpness as the signal is amplified. We then argue that under gradient descent with large enough step size, the sharpness will cross the stability threshold (in particular, at the parameters $b_o$), at which point the network will begin to *reintroduce* the noise variable. We complement this result with simulations demonstrating the proven behavior, as well as a comparison to almost identical behavior on a small MLP trained on a 5k subset of CIFAR-10.

We make several further experimental observations which seem relevant for interpreting the effect of these outliers and what that may imply more generally about NN optimization. We briefly list them here, with a more in-depth discussion in Appendix E: (i) sharpness often occurs overwhelmingly in the first few layers; (ii) batchnorm may smooth training, even if not the loss itself; (iii) for both GD and SGD, approximately half of training points go up in loss on each step; and (iv) different losses and label smoothing have predictable effects on sharpening.

## 4  The Interplay of Opposing Signals and Stochasticity

Full-batch GD is not used in practice when training NNs. It is therefore pertinent to ask what these findings imply about stochastic optimization. We begin by verifying that this pattern persists during SGD. Fig. 4 displays the losses for four opposing group pairs of a VGG-11-BN trained on CIFAR-10 with SGD batch size 128. We observe that the paired groups do exhibit clear opposing oscillatory patterns, but they do not alternate with every step, nor do they always move in opposite directions. This should not be surprising: we expect that not every batch will have a given signal in one direction or the other. For comparison, we include the *full* train loss in each figure—that is, including the points

not in the training batch. We see that the loss on the outliers has substantially larger variance; to confirm that this is not just because the groups have many fewer samples, we also plot the loss on a random subset of training points of the same size. We reproduce this plot with a VGG-11 without BN in Fig. 32 in the Appendix. Having verified that this oscillation on opposing signals still occurs in the stochastic setting, we conjecture that current best practices for neural network optimization owe at least some of their success to gracefully handling such imbalances. We investigate this possibility concretely for the Adam optimizer [21].



Figure 4: **VGG-11 trained with SGD on CIFAR-10.** Losses of paired outlier groups, along with the full train loss for comparison. The outliers' loss follow the same oscillatory pattern with large magnitude. See appendix for the same without BN.

**How Adam handles gradients with opposing signals.** To better understand their differences, Fig. 5 visualizes the parameter iterates of Adam and SGD with momentum on a ReLU MLP trained on a 5k subset of CIFAR-10, alongside those of GD and SGD. The top figure is the projection of these parameters onto the top eigenvector of the loss Hessian of the network trained with GD, evaluated at the first step where the sharpness crosses $2/\eta$. We observe that SGD tracks a similar path to GD, though adding momentum mitigates the oscillation somewhat. In contrast, the network optimized with Adam markedly departs from this pattern, smoothly oscillating along one side without crossing the middle. We identify several components of Adam which potentially contribute to this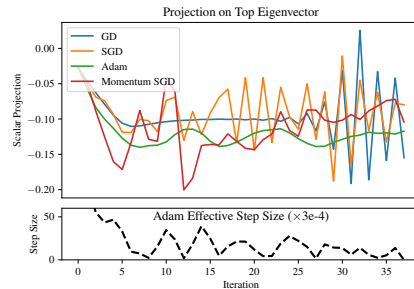 effect. For space constraints we list them here and expand upon each item in detail in Appendix G: (i) normalization means Adam takes very small steps along directions of large curvature, especially near the minimum; (ii) the "trust region" means there is not too much dependence on imbalanced opposing signals, and since it is not a descent method it isn't dominated by the direction "down the valley"; and (iii) Adam's gradient includes *dampening* in addition to the usual momentum term, which we argue is an important component.

To test whether our findings translate to practical gains, we design a variant of SGD which incorporates these insights. First, we use dampening $\tau = 0.9$ in addition to momentum. Second, we choose a global threshold: if the gradient magnitude for a given parameter is above this threshold, we take a fixed step size. Otherwise, we take a gradient step as normal. In Appendix H we show that this approach matches Adam when training a ResNet-56/110 on CIFAR-10 with learning rates across several orders of magnitude. We also compare the two methods for the initial phase of training of GPT-2 on the OpenWebText dataset—not only do they perform the same, their loss similarity suggests that their exact trajectory may be very similar. We find that the fraction of signed parameter steps is around 10% for the ResNet, and around 50% initially for the transformer, then gradually decaying.



Figure 5: **Projected iterates of an MLP on CIFAR-10. Top:** SGD closely tracks GD, bouncing across the valley; momentum somewhat mitigates the sharp jumps. Adam smoothly oscillates along one side. **Bottom:** Adam's effective step size drops sharply when moving too close or far from the valley floor.

## 5  Conclusion

The existence of outliers with such a significant yet non-obvious influence on neural network training raises as many questions as it answers. This work presents an initial investigation into their effect on various aspects of optimization, but there is still much more to understand. Though it is clear they have a large influence on training, less obvious is whether reducing their influence is *necessary* for improved optimization or simply coincides with it. At the same time, there is evidence that the behavior these outliers induce may serve as an important method of exploration and/or regularization. If so, another key question is whether these two effects can be decoupled—or if the incredible generalization ability of neural networks is somehow inherently tied to their instability.

# References

[1] Sanjeev Arora, Zhiyuan Li, and Abhishek Panigrahi. Understanding gradient descent on the edge of stability in deep learning. In *International Conference on Machine Learning*, pages 948–1024. PMLR, 2022.

[2] Boaz Barak, Benjamin L. Edelman, Surbhi Goel, Sham M. Kakade, eran malach, and Cyril Zhang. Hidden progress in deep learning: SGD learns parities near the computational limit. In *Advances in Neural Information Processing Systems*, 2022. URL https://openreview.net/forum?id=8XWP2ewX-im.

[3] Mikhail Belkin, Daniel Hsu, Siyuan Ma, and Soumik Mandal. Reconciling modern machine-learning practice and the classical bias–variance trade-off. *Proceedings of the National Academy of Sciences*, 116(32):15849–15854, 2019.

[4] Frederik Benzing. Gradient descent on neurons and its link to approximate second-order optimization. In *International Conference on Machine Learning*, pages 1817–1853. PMLR, 2022.

[5] Xiangning Chen, Chen Liang, Da Huang, Esteban Real, Kaiyuan Wang, Yao Liu, Hieu Pham, Xuanyi Dong, Thang Luong, Cho-Jui Hsieh, et al. Symbolic discovery of optimization algorithms. *arXiv preprint arXiv:2302.06675*, 2023.

[6] Jeremy Cohen, Simran Kaur, Yuanzhi Li, J Zico Kolter, and Ameet Talwalkar. Gradient descent on neural networks typically occurs at the edge of stability. In *International Conference on Learning Representations*, 2021. URL https://openreview.net/forum?id=jh-rTtvkGeM.

[7] Jeremy M Cohen, Behrooz Ghorbani, Shankar Krishnan, Naman Agarwal, Sourabh Medapati, Michal Badura, Daniel Suo, David Cardoze, Zachary Nado, George E Dahl, et al. Adaptive gradient methods at the edge of stability. *arXiv preprint arXiv:2207.14484*, 2022.

[8] Alex Damian, Eshaan Nichani, and Jason D. Lee. Self-stabilization: The implicit bias of gradient descent at the edge of stability. In *OPT 2022: Optimization for Machine Learning (NeurIPS 2022 Workshop)*, 2022. URL https://openreview.net/forum?id=enoU_Kp7Dz.

[9] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.

[10] Simon S Du, Wei Hu, and Jason D Lee. Algorithmic regularization in learning deep homogeneous models: Layers are automatically balanced. In *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018.

[11] Pierre Foret, Ariel Kleiner, Hossein Mobahi, and Behnam Neyshabur. Sharpness-aware minimization for efficiently improving generalization. In *International Conference on Learning Representations*, 2021. URL https://openreview.net/forum?id=6Tm1mposlrM.

[12] Stanislav Fort and Surya Ganguli. Emergent properties of the local geometry of neural loss landscapes. *arXiv preprint arXiv:1910.05929*, 2019.

[13] Matteo Gamba, Hossein Azizpour, and Mårten Björkman. On the lipschitz constant of deep networks and double descent. *arXiv preprint arXiv:2301.12309*, 2023.

[14] Behrooz Ghorbani, Shankar Krishnan, and Ying Xiao. An investigation into neural net optimization via hessian eigenvalue density. In *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*. PMLR, 09–15 Jun 2019. URL https://proceedings.mlr.press/v97/ghorbani19b.html.

[15] Minyoung Huh, Hossein Mobahi, Richard Zhang, Brian Cheung, Pulkit Agrawal, and Phillip Isola. The low-rank simplicity bias in deep networks. *arXiv preprint arXiv:2103.10427*, 2021.

[16] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International conference on machine learning*, pages 448–456. pmlr, 2015.

[17] Stanisław Jastrzębski, Zachary Kenton, Devansh Arpit, Nicolas Ballas, Asja Fischer, Yoshua Bengio, and Amos Storkey. Three factors influencing minima in sgd. *arXiv preprint arXiv:1711.04623*, 2017.

[18] Stanisław Jastrzębski, Zachary Kenton, Nicolas Ballas, Asja Fischer, Yoshua Bengio, and Amost Storkey. On the relation between the sharpest directions of DNN loss and the SGD step length. In *International Conference on Learning Representations*, 2019. URL `https://openreview.net/forum?id=SkgEaj05t7`.

[19] Stanisław Jastrzębski, Maciej Szymczak, Stanislav Fort, Devansh Arpit, Jacek Tabor, Kyunghyun Cho*, and Krzysztof Geras*. The break-even point on optimization trajectories of deep neural networks. In *International Conference on Learning Representations*, 2020. URL `https://openreview.net/forum?id=r1g87C4KwB`.

[20] Stanisław Jastrzębski, Devansh Arpit, Oliver Astrand, Giancarlo B Kerg, Huan Wang, Caiming Xiong, Richard Socher, Kyunghyun Cho, and Krzysztof J Geras. Catastrophic fisher explosion: Early phase fisher matrix impacts generalization. In *Proceedings of the 38th International Conference on Machine Learning*, Proceedings of Machine Learning Research. PMLR, 2021. URL `https://proceedings.mlr.press/v139/jastrzebski21a.html`.

[21] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.

[22] Dmitry Kopitkov and Vadim Indelman. Neural spectrum alignment: Empirical study. In *Artificial Neural Networks and Machine Learning–ICANN 2020: 29th International Conference on Artificial Neural Networks, Bratislava, Slovakia, September 15–18, 2020, Proceedings, Part II 29*, pages 168–179. Springer, 2020.

[23] Itai Kreisler, Mor Shpigel Nacson, Daniel Soudry, and Yair Carmon. Gradient descent monotonically decreases the sharpness of gradient flow solutions in scalar networks and beyond. In *Proceedings of the 40th International Conference on Machine Learning*, Proceedings of Machine Learning Research. PMLR, 23–29 Jul 2023. URL `https://proceedings.mlr.press/v202/kreisler23a.html`.

[24] Frederik Kunstner, Jacques Chen, Jonathan Wilder Lavington, and Mark Schmidt. Noise is not the main factor behind the gap between sgd and adam on transformers, but sign descent might be. In *The Eleventh International Conference on Learning Representations*, 2023. URL `https://openreview.net/forum?id=a65YK0cqH8g`.

[25] Aitor Lewkowycz, Yasaman Bahri, Ethan Dyer, Jascha Sohl-Dickstein, and Guy Gur-Ari. The large learning rate phase of deep learning: the catapult mechanism. *arXiv preprint arXiv:2003.02218*, 2020.

[26] Xinyan Li, Qilong Gu, Yingxue Zhou, Tiancong Chen, and Arindam Banerjee. Hessian based analysis of sgd for deep nets: Dynamics and generalization. In *Proceedings of the 2020 SIAM International Conference on Data Mining*, pages 190–198. SIAM, 2020.

[27] Yuanzhi Li, Colin Wei, and Tengyu Ma. Towards explaining the regularization effect of initial large learning rate in training neural networks. *Advances in Neural Information Processing Systems*, 32, 2019.

[28] Chao Ma and Lexing Ying. On linear stability of sgd and input-smoothness of neural networks. In *Advances in Neural Information Processing Systems*, volume 34, pages 16805–16817. Curran Associates, Inc., 2021. URL `https://proceedings.neurips.cc/paper_files/paper/2021/file/8c26d2fad09dc76f3ff36b6ea752b0e1-Paper.pdf`.

[29] Chao Ma, Daniel Kunin, Lei Wu, and Lexing Ying. Beyond the quadratic approximation: the multiscale structure of neural network loss landscapes. *arXiv preprint arXiv:2204.11326*, 2022.

[30] Lachlan Ewen MacDonald, Jack Valmadre, and Simon Lucey. On progressive sharpening, flat minima and generalisation. *arXiv preprint arXiv:2305.14683*, 2023.

[31] William Merrill, Nikolaos Tsilivis, and Aman Shukla. A tale of two circuits: Grokking as competition of sparse and dense subnetworks. In *ICLR 2023 Workshop on Mathematical and Empirical Understanding of Foundation Models*, 2023. URL `https://openreview.net/forum?id=8GZxtu46Kx`.

[32] Rotem Mulayoff and Tomer Michaeli. Unique properties of flat minima in deep networks. In *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*. PMLR, 13–18 Jul 2020. URL `https://proceedings.mlr.press/v119/mulayoff20a.html`.

[33] Preetum Nakkiran, Gal Kaplun, Dimitris Kalimeris, Tristan Yang, Benjamin L Edelman, Fred Zhang, and Boaz Barak. Sgd on neural networks learns functions of increasing complexity. *arXiv preprint arXiv:1905.11604*, 2019.

[34] Preetum Nakkiran, Gal Kaplun, Yamini Bansal, Tristan Yang, Boaz Barak, and Ilya Sutskever. Deep double descent: Where bigger models and more data hurt. In *International Conference on Learning Representations*, 2020. URL `https://openreview.net/forum?id=B1g5sA4twr`.

[35] Neel Nanda, Lawrence Chan, Tom Lieberum, Jess Smith, and Jacob Steinhardt. Progress measures for grokking via mechanistic interpretability. In *The Eleventh International Conference on Learning Representations*, 2023. URL `https://openreview.net/forum?id=9XFSbDPmdW`.

[36] Pascal Notsawo Jr, Hattie Zhou, Mohammad Pezeshki, Irina Rish, Guillaume Dumas, et al. Predicting grokking long before it happens: A look into the loss landscape of models which grok. *arXiv preprint arXiv:2306.13253*, 2023.

[37] Samet Oymak, Zalan Fabian, Mingchen Li, and Mahdi Soltanolkotabi. Generalization guarantees for neural networks via harnessing the low-rank structure of the jacobian. *arXiv preprint arXiv:1906.05392*, 2019.

[38] Yan Pan and Yuanzhi Li. Toward understanding why adam converges faster than sgd for transformers. *arXiv preprint arXiv:2306.00204*, 2023.

[39] Vardan Papyan. The full spectrum of deepnet hessians at scale: Dynamics with sgd training and sample size. *arXiv preprint arXiv:1811.07062*, 2018.

[40] Vardan Papyan. Measurements of three-level hierarchical structure in the outliers in the spectrum of deepnet hessians. In *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*. PMLR, 09–15 Jun 2019. URL `https://proceedings.mlr.press/v97/papyan19a.html`.

[41] Vardan Papyan. Traces of class/cross-class structure pervade deep learning spectra. *The Journal of Machine Learning Research*, 21(1):10197–10260, 2020.

[42] Alethea Power, Yuri Burda, Harri Edwards, Igor Babuschkin, and Vedant Misra. Grokking: Generalization beyond overfitting on small algorithmic datasets. *arXiv preprint arXiv:2201.02177*, 2022.

[43] Elan Rosenfeld, Pradeep Ravikumar, and Andrej Risteski. Domain-adjusted regression or: Erm may already learn features sufficient for out-of-distribution generalization. *arXiv preprint arXiv:2202.06856*, 2022.

[44] Levent Sagun, Leon Bottou, and Yann LeCun. Eigenvalues of the hessian in deep learning: Singularity and beyond. *arXiv preprint arXiv:1611.07476*, 2016.

[45] Levent Sagun, Utku Evci, V Ugur Guney, Yann Dauphin, and Leon Bottou. Empirical analysis of the hessian of over-parametrized neural networks. *arXiv preprint arXiv:1706.04454*, 2017.

[46] Shibani Santurkar, Dimitris Tsipras, Andrew Ilyas, and Aleksander Madry. How does batch normalization help optimization? *Advances in neural information processing systems*, 31, 2018.

[47] Andrew M Saxe, James L McClelland, and Surya Ganguli. Exact solutions to the nonlinear dynamics of learning in deep linear neural networks. *arXiv preprint arXiv:1312.6120*, 2013.

[48] Harshay Shah, Kaustav Tamuly, Aditi Raghunathan, Prateek Jain, and Praneeth Netrapalli. The pitfalls of simplicity bias in neural networks. *Advances in Neural Information Processing Systems*, 33:9573–9585, 2020.

[49] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.

[50] Vimal Thilak, Etai Littwin, Shuangfei Zhai, Omid Saremi, Roni Paiss, and Joshua Susskind. The slingshot mechanism: An empirical study of adaptive optimizers and the grokking phenomenon. *arXiv preprint arXiv:2206.04817*, 2022.

[51] Guillermo Valle-Perez, Chico Q. Camargo, and Ard A. Louis. Deep learning generalizes because the parameter-function map is biased towards simple functions. In *International Conference on Learning Representations*, 2019. URL `https://openreview.net/forum?id=rye4g3AqFm`.

[52] Shengjie Wang, Abdel-rahman Mohamed, Rich Caruana, Jeff Bilmes, Matthai Plilipose, Matthew Richardson, Krzysztof Geras, Gregor Urban, and Ozlem Aslan. Analysis of deep neural networks with extended data jacobian matrix. In *International Conference on Machine Learning*, pages 718–726. PMLR, 2016.

[53] Zixuan Wang, Zhouzi Li, and Jian Li. Analyzing sharpness along GD trajectory: Progressive sharpening and edge of stability. In *Advances in Neural Information Processing Systems*, 2022. URL `https://openreview.net/forum?id=thgItcQrJ4y`.

[54] Lei Wu, Chao Ma, and Weinan E. How sgd selects the global minima in over-parameterized learning: A dynamical stability perspective. In *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018. URL `https://proceedings.neurips.cc/paper_files/paper/2018/file/6651526b6fb8f29a00507de6a49ce30f-Paper.pdf`.

[55] Lei Wu, Mingze Wang, and Weijie J Su. The alignment property of SGD noise and how it helps select flat minima: A stability analysis. In *Advances in Neural Information Processing Systems*, 2022. URL `https://openreview.net/forum?id=rUc8peDIM45`.

[56] Chen Xing, Devansh Arpit, Christos Tsirigotis, and Yoshua Bengio. A walk with sgd. *arXiv preprint arXiv:1802.08770*, 2018.

[57] Jingzhao Zhang, Sai Praneeth Karimireddy, Andreas Veit, Seungyeon Kim, Sashank Reddi, Sanjiv Kumar, and Suvrit Sra. Why are adaptive methods good for attention models? *Advances in Neural Information Processing Systems*, 33:15383–15393, 2020.

[58] Xingyu Zhu, Zixuan Wang, Xiang Wang, Mo Zhou, and Rong Ge. Understanding edge-of-stability training dynamics with a minimalist example. In *The Eleventh International Conference on Learning Representations*, 2023. URL `https://openreview.net/forum?id=p7EagBsMAEO`.

[59] Zhanxing Zhu, Jingfeng Wu, Bing Yu, Lei Wu, and Jinwen Ma. The anisotropic noise in stochastic gradient descent: Its behavior of escaping from sharp minima and regularization effects. In *Proceedings of the 36th International Conference on Machine Learning*, Proceedings of Machine Learning Research, pages 7654–7663. PMLR, 09–15 Jun 2019. URL `https://proceedings.mlr.press/v97/zhu19e.html`.

# A  Discussion and Future Work

Many of the observations we make in this paper are not new, having been described in various prior works. Rather, this work identifies a possible *higher-order cause* which neatly ties these findings together. There are also many works which pursue a more theoretical understanding of each of these phenomena independently. Such analyses begin with a set of assumptions (on the data, in particular) and prove that the given behavior follows. In contrast, this work *begins* by identifying a condition—the presence of opposing signals—which we argue is likely a major cause of these behaviors. These two are not at odds: we believe in many cases our result serves as direct evidence for the validity of these modeling assumptions and that it may enable even more fine-grained analyses. This work provides an initial investigation which we hope will inspire future efforts towards a more complete understanding.

We now highlight some connections to these earlier findings. More general related work can be found in the next section.

**Heavy-tailed loss spectrum.**   Earlier studies of the loss landscape noted a small group of very large outlier Hessian eigenvalues or Jacobian singular values (e.g. Sagun et al. 44, 45, Papyan 39, see Appendix B for more). Our method of identifying these paired groups, along with the metrics tracked in Fig. 29, indicate that these outlier directions in the spectrum are precisely the directions with opposing signals in the gradient and that this pattern may be key to better understanding the generalization ability of NNs trained with SGD.

**Progressive sharpening and the edge of stability.**   More recent focus has shifted to the top Hessian eigenvalue(s), where it was empirically observed that their magnitude (the loss "sharpness") grows during training [18, 19, 6] (so-called *progressive sharpening*), leading to rapid oscillation in weight space [56, 18]. Cohen et al. [6] also found that for GD this coincides with a consistent yet non-monotonic decrease in training loss over long timescales, which they named the *edge of stability*. We observe that prior analyses have proven the *occurrence* of progressive sharpening and the edge of stability under various assumptions [1, 53], but the underlying *cause* has not been made clear. Our discussion, experiments, and theoretical analysis in Section 3 provide strong evidence for a genuine cause which aligns with several of these existing modeling assumptions. Roughly, it seems that progressive sharpening occurs when the network learns to rely on (or *not* rely on) opposing signals in a very specific way, while simultaneously amplifying other, more useful signal. This growth in sensitivity means a small parameter change modifying how opposing signals are used can massively increase loss. This leads to intermittent instability orthogonal to the "valley floor", accompanied by gradual training loss decay and occasional spikes as described by the toy example in Fig. 2 and depicted on real data in Fig. 1. Empirically, this oscillation seems somewhat independent of movement parallel to the floor (see Appendix H), but further study of the precise dynamics is needed.

**Spurious correlations, grokking, and slingshotting.**   In images, the features corresponding to opposing signals match the traditional picture of "spurious correlations" surprisingly closely—it could be that a network maintaining balance or diverging along a direction also determines whether it continues to use a "spurious" feature or is forced to find an alternative way to minimize loss. Indeed, the exact phenomenon of a network "slingshotting" to a new region with improved generalization has been directly observed [54, 25, 20, 50]. *Grokking* [42], whereby a network learns to generalize long after memorizing the training set, is closely related. Several works have shown that grokking is a "hidden" phenomenon, with gradual amplification of generalizing subnetworks [2, 35, 31]; it has even been noted to co-occur with weight oscillation [36]. Our experiments in Section 4 and Appendix H show that the influence of opposing signals obscures the behavior of the rest of the network, offering one possible explanation.

**Simplicity bias and double descent.**   Nakkiran et al. [33] observed that NNs learn functions of greater complexity throughout training. Our experiments—particularly the slow decay in the norm of the feature embedding of opposing signals—lead us to believe it would be more correct to say that they *unlearn* simple functions, which enables more complex subnetworks with smaller magnitude and better performance to take over. At first this seems at odds with the notion of *simplicity bias* [51, 48], defined broadly as a tendency of networks to rely on simple functions of their inputs. However, it does seem to be the case that the network will use the simplest (e.g., largest norm)

features that it can, so long as such features allow it to approach zero training loss; otherwise it may eventually diverge. This tendency also suggests a possible explanation for *double descent* [3, 34]: even after interpolation, the network pushes towards greater confidence and the weight layers continue to balance [47, 10], increasing sharpness. This could lead to oscillation, pushing the network to learn new features which generalize better [54, 27, 43, 50]. This behavior would also be more pronounced for larger networks because they exhibit greater sharpening. Note that the true explanation is not quite so straightforward: generalization is sometimes improved via methods that *reduce* oscillation (like loss smoothing), implying that this behavior is not always advantageous. A better understanding of these nuances is an important subject for future study.

**Sharpness-Aware Minimization**   Another connection we think merits further inquiry is Sharpness-Aware Minimization (SAM) [11], which is known to improve generalization of neural networks for reasons still not fully understood. In particular, the better-performing variant is 1-SAM, which takes positive gradient steps on each training point in the batch individually. It it evident that several of these updates will point along directions of steepest descent/ascent orthogonal to the valley floor (and, if not normalized, the updates may be *very* large). Thus it may be that 1-SAM is in some sense "simulating" divergence out of this valley in both directions, enabling exploration in a manner that would not normally be possible until the sharpness grows large enough. In contrast, standard SAM would only take this step in one of the two directions, or perhaps not at all if the opposing signals are equally balanced in the minibatch. SAM's intermediate steps would also encourage the network to downweight these features faster. These possibilities seem a promising direction for further exploration.

# B   Related Work

**Characterizing the NN loss landscape.**   Earlier studies of the loss landscape commonly identified a heavy-tailedness with a small group of very large outlier Hessian eigenvalues or Jacobian singular values [44, 45, 39, 37, 40, 12, 14, 26, 41, 22]. Later efforts focused on concretely linking these observations to corresponding behavior, often with an emphasis on SGD's bias towards particular solutions [54, 17, 19] and what this may imply about its resulting generalization [18, 59, 55]. Our method for identifying these paired groups, along with Fig. 29, indicates that these outlier directions in the Hessian/Jacobian spectrum are precisely the directions with opposing signals in the gradient, and that this pattern may be key to better understanding the generalization ability of NNs trained with SGD.

**Progressive sharpening and the edge of stability.**   Shifting away from the overall structure, more recent focus has been specifically on top eigenvalue(s), where it was empirically observed that their magnitude (the loss "sharpness") grows when training with SGD [18, 19] and GD [22, 6] (so-called "progressive sharpening"). This leads to rapid oscillation in weight space [56, 18, 6, 7]. Cohen et al. [6] also found that for GD this coincides with a consistent yet non-monotonic decrease in training loss over long timescales, which they named the "edge of stability"; moreover, they noted that this behavior runs contrary to our traditional understanding of NN convergence. Many works have since investigated the possible origins of this phenomenon [58, 23]. Several of these are deeply related to our findings: Ma et al. [29] connect this behavior to the existence of multiple "scales" of losses; the outliers we identify corroborate this point. Damian et al. [8] prove that GD implicitly regularizes the sharpness—we identify a conceptually distinct source of such regularization, as described in Section 3. Arora et al. [1] show under some conditions that the GD trajectory follows a minimum-loss manifold towards lower curvature regions. This is consistent with our findings, and we believe this manifold to be precisely the path which evenly balances the opposing gradients. Wang et al. [53] provide another thorough analysis of NN training dynamics at the edge of stability; their demonstrated phases closely align with our own. They further observe that this sharpening coincides with a growth in the norm of the last layer, which was also noted by MacDonald et al. [30]. Our proposed explanation for the effect of opposing signals offers some insight into this relationship, but further investigation is needed.
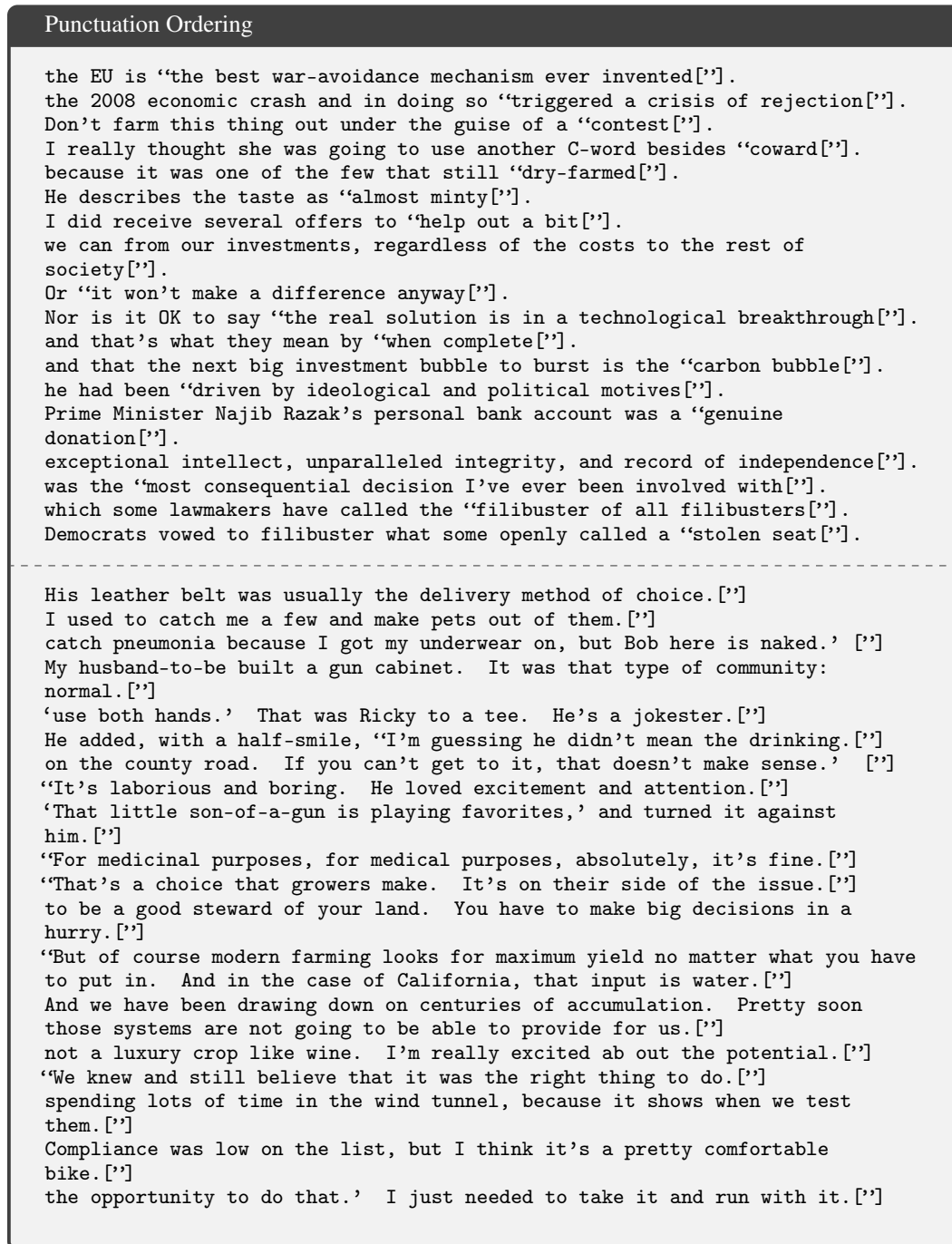
# C    Examples of Opposing Signals in Text

```
┌─────────────────────────────────────────────────────────────────────────────┐
│ Punctuation Ordering                                                          │
├─────────────────────────────────────────────────────────────────────────────┤
│                                                                               │
│  the EU is "the best war-avoidance mechanism ever invented["].                │
│  the 2008 economic crash and in doing so "triggered a crisis of rejection["]. │
│  Don't farm this thing out under the guise of a "contest["].                  │
│  I really thought she was going to use another C-word besides "coward["].     │
│  because it was one of the few that still "dry-farmed["].                      │
│  He describes the taste as "almost minty["].                                  │
│  I did receive several offers to "help out a bit["].                          │
│  we can from our investments, regardless of the costs to the rest of          │
│  society["].                                                                  │
│  Or "it won't make a difference anyway["].                                    │
│  Nor is it OK to say "the real solution is in a technological breakthrough["].│
│  and that's what they mean by "when complete["].                              │
│  and that the next big investment bubble to burst is the "carbon bubble["].   │
│  he had been "driven by ideological and political motives["].                 │
│  Prime Minister Najib Razak's personal bank account was a "genuine            │
│  donation["].                                                                 │
│  exceptional intellect, unparalleled integrity, and record of independence["].│
│  was the "most consequential decision I've ever been involved with["].        │
│  which some lawmakers have called the "filibuster of all filibusters["].      │
│  Democrats vowed to filibuster what some openly called a "stolen seat["].     │
│- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -│
│  His leather belt was usually the delivery method of choice.["]               │
│  I used to catch me a few and make pets out of them.["]                       │
│  catch pneumonia because I got my underwear on, but Bob here is naked.' ["]   │
│  My husband-to-be built a gun cabinet.  It was that type of community:         │
│  normal.["]                                                                   │
│  'use both hands.'  That was Ricky to a tee.  He's a jokester.["]             │
│  He added, with a half-smile, "I'm guessing he didn't mean the drinking.["]   │
│  on the county road.  If you can't get to it, that doesn't make sense.'  ["]  │
│  "It's laborious and boring.  He loved excitement and attention.["]           │
│  'That little son-of-a-gun is playing favorites,' and turned it against       │
│  him.["]                                                                      │
│  "For medicinal purposes, for medical purposes, absolutely, it's fine.["]     │
│  "That's a choice that growers make.  It's on their side of the issue.["]     │
│  to be a good steward of your land.  You have to make big decisions in a       │
│  hurry.["]                                                                    │
│  "But of course modern farming looks for maximum yield no matter what you have│
│  to put in.  And in the case of California, that input is water.["]           │
│  And we have been drawing down on centuries of accumulation.  Pretty soon      │
│  those systems are not going to be able to provide for us.["]                 │
│  not a luxury crop like wine.  I'm really excited ab out the potential.["]    │
│  "We knew and still believe that it was the right thing to do.["]             │
│  spending lots of time in the wind tunnel, because it shows when we test       │
│  them.["]                                                                     │
│  Compliance was low on the list, but I think it's a pretty comfortable         │
│  bike.["]                                                                     │
│  the opportunity to do that.'  I just needed to take it and run with it.["]   │
│                                                                               │
└─────────────────────────────────────────────────────────────────────────────┘
```

Figure 6: **Examples of opposing signals in text.** Found by training GPT-2 on a subset of Open-WebText. Sequences are on separate lines, the token in brackets is the target and all prior tokens are (the end of the) context. As both standards are used, it is not always clear whether punctuation will come before or after the end of a quotation (we include the period after the quote for clarity—the model does not condition on it). Note that the double quotation is encoded as the *pair* of tokens [447, 251], and the loss oscillation is occurring for sequences that end with this pair, either before (top) or after (bottom) the occurrence of the period token (13).

```
New Line or 'the' After Colon

In order to prepare your data, there are three things to do:[\n]
in the FP lib of your choice, namely Scalaz or Cats. It looks like this:[\n]
Let the compiler guide you, it will only accept one implementation:[\n]
Salcedo said of the work:[\n]
Enter your email address:[\n]
According to the CBO update:[\n]
Here's how the Giants can still make the playoffs:[\n]
described how he copes with his condition in an interview with The
Telegraph:[\n]
Here's a list of 5 reasons as to why self diagnosis is valid:[\n]
successive Lambda invocations. It looks more or less like this:[\n]
data, there are three things to do:[\n]
4.2 percent in early 2018.\n\nAccording to the CBO update:[\n]
other than me being myself.''\n\nWATCH:[\n]
is to make the entire construction plural.\n\nTwo recent examples:[\n]
We offer the following talking points to anyone who is attending the
meeting:[\n]
is on the chopping block - and at the worst possible moment:[\n]
- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -
as will our MPs in Westminster. But to me it is obvious:  [the]
The wheelset is the same as that on the model above:  [the]
not get so engrained or in a rut with what I had been doing.  Not to worry:
[the]
polemics against religion return in various ways to one core issue:  [the]
which undergirds all other acts of love, both divine and human:  [the]
integrate fighters from the Kurds' two main political parties:  [the]
robs this incredible title of precisely what makes it so wonderful:  [the]
you no doubt noticed something was missing:  [the]
Neil Gorsuch's 'sexist' comments on maternity leave:  [the]
```

Figure 7: **Examples of opposing signals in text.** Found by training GPT-2 on a subset of OpenWeb-Text. Sequences are on separate lines, the token in brackets is the target and all prior tokens are (the end of the) context. Sometimes a colon occurs mid-sentence—and is often followed by "the"—other times it announces the start of a new line. The model must *unlearn* ": $\mapsto$ [\n]" versus ": $\mapsto$ [the]" and instead use other contextual information.



Figure 8: Loss of GPT-2 on the above opposing signals.

14

# D    Reproducing Fig. 3 in Other Settings

Though colors are straightforward, for some opposing signals such as grass texture it is not clear how to produce a synthetic image which properly captures what precisely the model is latching on to. Instead, we identify a real image which has as much grass and as little else as possible, with the understanding that the additional signal in the image could affect the results. We depict the grass image alongside the plots it produced.

## D.1    ResNet-18 Trained with GD on Other Inputs



Figure 9: ResNet-18 on a red color block.



Figure 10: ResNet-18 on a green color block. As this color seems unnatural, we've included two examples of relevant images in the dataset.



Figure 11: Examples of images with the above green color.



Figure 12: ResNet-18 on a white color block.

**Input Image**



Figure 13: ResNet-18 on a black color block.



Figure 14: ResNet-18 on an image with mostly grass texture.

## D.2 VGG-11-BN Trained with GD

For VGG-11, we found that the feature norm of the embedded images did not decay nearly as much over the course of training. We expect this has to do with the lack of a residual component. However, for the most part these features do still follow the pattern of a rapid increase, followed by a marked decline.



Figure 15: VGG-11-BN on a sky color block.



Figure 16: VGG-11-BN on a red color block.



Figure 17: VGG-11-BN on a green color block. See above for two examples of relevant images in the dataset.



Figure 18: VGG-11-BN on an image with mostly grass texture.

17

Figure 19: VGG-11-BN on a white color block.



Figure 20: VGG-11-BN on a black color block.

## D.3 VGG-11-BN with Small Learning Rate

Here we see that oscillation at the edge of stability is an important regularizer to prevent the network from continuously upweighting opposing signals. As described in the main body, stepping too far in one direction causes an imbalanced gradient between the two opposing signals. Since the group which now has a larger loss is also the one which suffers from the use of the feature, the network is encouraged to downweight its influence. If we use a small enough learning rate that optimization closely tracks gradient flow, this regularization does not occur and the feature norms grow continuously, leading to over-reliance on this features and poor generalization.



Figure 21: VGG-11-BN on a sky color block with learning rate 0.008, approximating gradient flow.



Figure 22: VGG-11-BN on a red color block with learning rate 0.008, approximating gradient flow.



Figure 23: VGG-11-BN on a green color block with learning rate 0.0008, approximating gradient flow.

## D.4 ResNet-18 Trained with Full-Batch Adam

Finally, we plot the same figures for a ResNet-18 trained with full-batch Adam. We see that Adam consistently and quickly reduces the norm of these features, especially for more complex features such as texture (note when comparing to plots above that the maximum iteration differs).

Figure 24: VGG-11-BN on an image with mostly grass texture with learning rate 0.0008, approximating gradient flow.



Figure 25: ResNet-18 on a sky color block trained with Adam.



Figure 26: ResNet-18 on a red color block trained with Adam.



Figure 27: ResNet-18 on a green color block trained with Adam.



Figure 28: ResNet-18 on an image with mostly grass texture trained with Adam.

# E  Discussion of Additional Experimental Findings

**Implications for the effect of loss smoothing or other losses, as well as various activations.**    We found that the sharpness in ResNets occurred overwhelmingly in the first convolutional layer, after the first few training steps where it was in the last layer. For VGG, it occurred mostly in the first few layers. In transformers, curvature typically was most concentrated in the initial embedding layer and the first few MLP projection layers. Generally, it seems that the components of the network which *interact most directly with the input* have the most significant sharpness—particularly if they also perform dimensionality reduction. This is consistent with our understanding of what causes said sharpness.

**Batchnorm may smooth training, even if not the loss itself.** We also found that adding an additional batchnorm (BN) layer [16] before the last ResNet layer reduced its initial sharpness in that location. Cohen et al. [6] noted that BN does not prevent networks from reaching the edge of stability and concluded, contrary to [46], that BN does not smooth the loss landscape; we **conjecture** that the effect of BN depends on the use of GD vs. SGD. Specifically, our findings hint at a possible benefit of BN which applies *only* to minibatches: reducing the influence of imbalanced opposing signals. This suggests that in fact BN may smooth the *optimization trajectory* of neural networks, rather than the loss itself (this is consistent with the distinction made by Cohen et al. [6] between regularity and smoothness). In Section 4 we demonstrate that Adam also smooths the optimization trajectory and that minor changes to emulate this effect can aid stochastic optimization.

**For both GD and SGD, approximately half of training points go up in loss on each step.**    Though only the outliers are wildly oscillating, many more images contain some small component of the features they exemplify. Fig. 31 tracks the fraction of points which increase in loss on each step—to some extent, a small degree of oscillation appears to be happening to the entire dataset.

**Different losses and label smoothing have predictable effects on sharpening.**    [30] note that label smoothing the cross-entropy loss reduces sharpening. This is reasonably explained by the fact that smoothing reduces the loss suffered by extreme overconfidence, and therefore the loss for a given opposing signal will not be as sharp. This also hints at why logistic loss may be more suitable for NN optimization, because it only has substantial curvature around $x = 0$ so, unlike square or exponential loss, large steps will not massively increase sharpness. We expect a similar property may contribute to the relative performance of various activations (e.g., ReLU).

# F  Additional Figures



Figure 29: **Tracking other metrics which characterize outliers with opposing signals.** Maximal per-step change in loss relates to other useful metrics, such as per-sample gradient norm and curvature. We combine each pair of groups in Fig. 1 to create training subsets which each exemplify one "signal": we see that these samples are also significant outliers according to the other metrics. (For a point $x$, "Curvature on Top Full Loss Eigenvector" is defined as $v^\top H(x)v$, where $v$ is the top eigenvector of the full loss Hessian and $H(x)$ is the Hessian of the loss on $x$ alone.)

## F.1  Training a Small MLP on a Chebyshev Polynomial

Following Cohen et al. [6] we train a small MLP to fit a Chebyshev polynomial on evenly spaced points in the interval $[-1, 1]$ (Fig. 30 in Appendix). This data has no "outliers" in the traditional

sense, and it is not immediately clear what opposing signals are present. Nevertheless, we observe the same alternating behavior: we find a pair where one group is a small interval of $x$-values and the opposing group contains its neighbors, all in the range $[-1, -0.5]$. This suggests that the network has internal activations which are heavily influential only for more negative $x$-values. In this context, these two groups are the outliers.
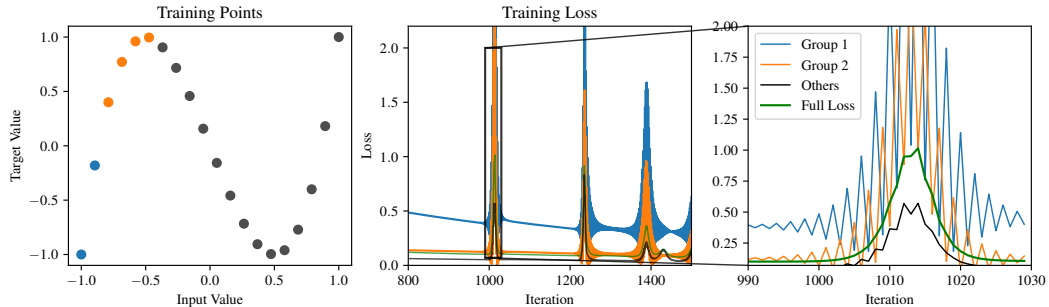


Figure 30: **Opposing signals when fitting a Chebyshev polynomial with a small MLP.** Though the data lacks traditional "outliers", it is apparent that the network has some features which are most influential only on the more negative inputs (or whose effect is otherwise cancelled out by other features). Since the correct use of this feature is opposite for these two groups, they provide opposing signals.



Figure 31: The fraction of overall training points which increase in loss on any given step. For both SGD and GD, it hovers around 0.5 (VGG without batchnorm takes a long time to reach the edge of stability). Though only the outliers are wildly swinging in loss, many more images contain *some* small component of the features they exemplify, and so these points also oscillate in loss at a smaller scale.

Figure 32: We reproduce Fig. 4 without batch normalization.



(a) A 3-layer ReLU MLP trained on a 5k-subset of CIFAR-10.

(b) Our model: a 2-layer linear network trained on mostly Gaussian data with opposing signals.

Figure 33: We compare a small ReLU MLP on a subset of CIFAR-10 to our simple model of linear regression with a two-layer network.

# G  Discussion on the Advantages of Adam

**Advantage 1: Smaller steps along high curvature directions.**  Adam's normalization causes smaller steps along the top eigenvector, especially near the minimum. The lower plot in Fig. 5 shows that the effective step size in this direction—i.e., the absolute inner product of the parameter-wise step sizes and the top eigenvector—rapidly drops to zero as the iterates approach the valley floor (in the opposite direction, the gradient negates the momentum for the same effect). We conjecture that general normalization may not be essential to Adam's performance; we even expect it could be somewhat harmful by limiting exploration. On the other hand, normalizing steps by curvature *parameter-wise* does seem important; Pan and Li [38] argue the same and show that parameter-wise gradient clipping improves SGD substantially. We highlight why this may be useful in the next point.

**Advantage 2: Managing heavy-tailed gradients and avoiding steepest descent.**  Zhang et al. [57] identified the "trust region" as an important contributor to Adam's success in attention models, pointing to heavy-tailed noise in the stochastic gradients. More recently, Kunstner et al. [24] argued that Adam's superiority does not come from better handling noise, which they supported by experimenting with large batch sizes. Our result reconciles these contradictory claims by showing that **the difficulty is not heavy-tailed *noise*, but strong, directed (and perhaps imbalanced) opposing signals.** Unlike traditional "gradient noise", larger batch sizes may not reduce the effect of these signals—that is, the gradient is heavy-tailed (across parameters) even without being stochastic. We also point to a related property of Adam: the largest steps emulate Sign SGD, which is notably *not* a descent method. Fig. 5 shows that Adam's steps are more parallel to the valley floor than those of steepest descent. **Thus it seems advantageous to *intentionally* avoid steps along the gradient which point towards the local minimum**—though not necessary so long as the step size is not too small. Indeed, Benzing [4] observe that true second order methods perform worse than SGD on NNs, and Kunstner et al. [24] show that Adam shares some behavior with Sign SGD with momentum. This point is also consistent with the observed generalization benefits of a large learning rate for SGD on NNs [19]; in fact, opposing signals naturally fit the concept of "easy-to-generalize" features as modeled by Li et al. [27]. In Appendix D.4 we again track the logits and feature norms of different color blocks on a ResNet-18, this time trained with full-batch Adam. We see that it more quickly and consistently reduces the norm of these features— further evidence that intentionally avoiding steepest descent prevents over-reliance on these features.

**Advantage 3: Dampening.**  Lastly, Adam's third important factor: downweighting the most recent gradient. Traditional SGD with momentum $\beta < 1$ takes a step which weights the current gradient by $\frac{1}{1+\beta} > \frac{1}{2}$. Though this makes intuitive sense, our results imply that heavily weighting the most recent gradient can be problematic. Instead, we expect an important addition is *dampening*, which multiplies the stochastic gradient at each step by some $(1 - \tau) < 1$. We observe that Adam's (unnormalized) gradient is equivalent to SGD with momentum and dampening both equal to $\beta_1$, plus a debiasing step. Recently proposed alternatives also include dampening in their momentum update but do not explicitly identify the distinction [57, 38, 5].

# H   Comparing our Variant of SGD to Adam

---

**Algorithm 1** SplitSGD

> **input:** Initial parameters $\theta_0$, SGD step size $\eta_1$, SignSGD step size $\eta_2$, momentum $\beta$, dampening $\tau$, threshold $r$.
> **initialize:** $m_0 = \mathbf{0}$.
> **for** $t \leftarrow 1, \ldots, T$ **do**
> $\quad g_t \leftarrow \nabla_\theta L_t(\theta_{t-1})$ $\hfill \triangleright$ Get stochastic gradient
> $\quad m_t \leftarrow \beta m_{t-1} + (1-\tau)g_t$ $\hfill \triangleright$ Update momentum with dampening
> $\quad \hat{m}_t \leftarrow m_t/(1-\tau^t)$ $\hfill \triangleright$ Debias
> $\quad v_{\text{mask}} \leftarrow \mathbf{1}\{|\hat{m}_t| \leq r\}$ $\hfill \triangleright$ Split parameters by threshold
> $\quad \theta_t \leftarrow \theta_{t-1} - \eta_1(\hat{m}_t \odot v_{\text{mask}}) - \eta_2(\text{sign}(\hat{m}_t) \odot (1 - v_{\text{mask}}))$ $\hfill \triangleright$ unmasked SGD, masked SignSGD
> **end for**

---

As described in the main text, we find that simply including dampening and taking a fixed step size on gradients above a certain threshold results in performance matching that of Adam for the experiments we tried. We found that setting this threshold equal to the $q = .1$ quantile of the first gradient worked quite well—this was about `1e-4` for the ResNet-56/110 and `1e-6` for GPT-2.

Simply to have something to label it with, we name the method SplitSGD, because it performs SGD and SignSGD on different partitions of the parameters. The precise method is given above in Algorithm 1. We reiterate that we are not trying to suggest a new method—our goal is only to demonstrate the insight gained from knowledge of opposing signals' influence on NN optimization. For all plots, $\beta$ represents the momentum parameter and $\tau$ is dampening. Adam has a single parameter $\beta_1$ which represents both simultaneously, which we fix at 0.9, and we do the same for SplitSGD. As in Algorithm 1, we let $\eta_1$ refer to the learning rate for standard SGD on the parameters with gradient below the magnitude threshold, and $\eta_2$ to the learning rate for the remainder which are optimized with SignSGD.

## H.1   SplitSGD on ResNet



Figure 34: Standard SGD with varying learning rates and momentum/dampening parameters on a ResNet-56 on CIFAR-10, with one run of SplitSGD for comparison. Omitted SGD hyperparameter combinations performed much worse. Notice that SGD is extremely sensitive to hyperparameters. Rightmost plot is the fraction of parameters with fixed step size by SplitSGD.

We begin with a comparison on ResNets trained on CIFAR-10. Fig. 34 compares SplitSGD to standard versions of SGD with varying momentum and dampening on a ResNet-56. As expected, SGD is extremely sensitive to hyperparameters, particularly the learning rate, and even the best choice in a grid search underperforms SplitSGD. Furthermore, the rightmost plot depicts the fraction of parameters for which SplitSGD takes a fixed-size signed step. This means that after the first few training steps, 70-80% of the parameters are being optimized simply with standard SGD (with $\beta = \tau = 0.9$).

Next, Fig. 35 plots SplitSGD with varying $\eta_1$ and $\eta_2$ fixed at .001. This is compared to Adam with learning rate .005, which was chosen via oracle grid search. Even though the SGD learning rate
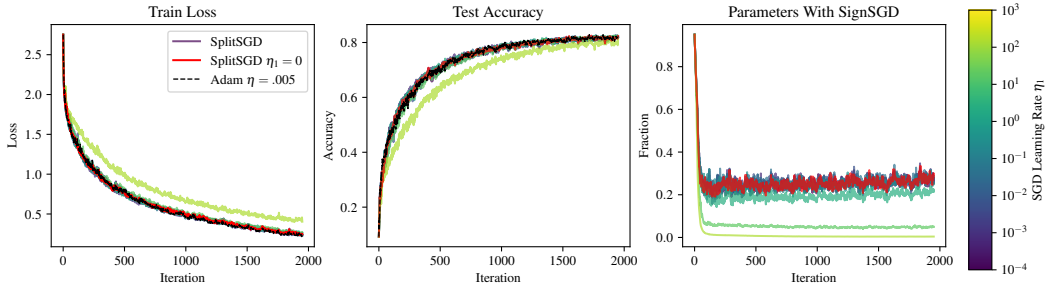
Figure 35: SplitSGD with varying SGD learning rates $\eta_1$ versus Adam on a ResNet-56 on CIFAR-10. The SignSGD learning rate is fixed at $\eta_2 = .001$; Adam uses $\eta = .005$, which was found to be the best performing choice via oracle selection grid search. The rightmost plot is the fraction of parameters with fixed step size by SplitSGD—that is, 1 minus this value is the fraction of parameters taking a regular gradient step with step size as given in the legend. This learning rate ranges over several orders of magnitude, is used for ~70-80% of parameters, and can even be set to 0, with no discernible difference in performance.

$\eta_1$ ranges over *seven orders of magnitude* and is used for ~70-80% of parameters, we see no real difference in the train loss or test accuracy of SplitSGD. In fact, we find that we can even eliminate it completely! This suggests that for most of parameters and most of training, it is only a small fraction of parameters in the entire network which are influencing the overall performance. We posit a deeper connection here to the "hidden" progress described in grokking [2, 35]—if the correct subnetwork and its influence on the output grows slowly during training, that behavior will not be noticeable until the dominating signals are first downweighted.

Figs. 36 and 37 depict the train loss and test accuracy of Adam and SplitSGD for varying learning rates (the standard SGD learning rate $\eta_1$ is fixed at 0.1). We see that SplitSGD is at least as robust as Adam to learning rate choice, if not more. The results also suggest that SplitSGD benefits from a slightly smaller learning rate than Adam, which we attribute to the fact that it will *always* take step sizes of that fixed size, whereas the learning rate for Adam represents an upper bound on the step size for each parameter.



Figure 36: Train loss of Adam and SplitSGD for varying learning rates. The regular SGD step size for SplitSGD is fixed at 0.1. SplitSGD seems at least as robust to choice of learning rate as Adam, and it appears to benefit from a slightly smaller learning rate because it cannot adjust per-parameter.



Figure 37: Test accuracy of Adam and SplitSGD for varying learning rates. The regular SGD step size for SplitSGD is fixed at 0.1. SplitSGD seems at least as robust to choice of learning rate as Adam, and it appears to benefit from a slightly smaller learning rate because it cannot adjust per-parameter.

We repeat these experiments with a ResNet-110, with similar findings. Fig. 38(a) compares the train loss and test accuracy of SGD with $\beta = 0.9, \tau = 0$ to Adam, and again the sensitivity of this optimizer to learning rate is clear. Fig. 38(b) compares Adam to SplitSGD (both with fixed-step learning rate .0003) but ablates the use of dampening: we find that the fixed-size signed steps appear to be more important for early in training, while dampening is helpful for maintaining performance later. It is not immediately clear what causes this bifurcation, nor if it will necessarily transfer to attention models.

Finally, Fig. 39(a) compares Adam to the full version of SplitSGD; we see essentially the same performance, and furthermore SplitSGD maintains its robustness to the choice of standard SGD learning rate.



(a) Adam versus standard SGD with Momentum. SGD remains extremely sensitive to choice of learning rate.

(b) Adam vs. SplitSGD with $\tau = 0$. Fixed-size learning rate for both is .0003.



(a) Adam vs. SplitSGD with $\tau = 0.9$. Fixed-size learning rate for both is .0003.

(b) The fraction of parameters for which a fixed-size signed step was taken for each gradient step.

## H.2  SplitSGD on GPT-2

For the transformer, we use the public nanoGPT repository which trains GPT-2 on the OpenWebText dataset. As a full training run would be too expensive, we compare only for the early stage of optimization. All hyperparameters are the defaults from that repository, with the SGD learning rate $\eta_1$ set equal to the other learning rate $\eta_2$. We observe that not only do the two methods track each other closely in training loss, it appears that they experience *exactly* the same oscillations. Though we do not track the parameters themselves, this suggests that these two methods follow very similar optimization trajectories as well, which we believe is an intriguing possibility worth further study.

Figure 40: Adam versus SplitSGD on the initial stage of training GPT-2 on the OpenWebText dataset, and the fraction of parameters with a fixed-size signed step. All hyperparameters are the defaults from the nanoGPT repository. Observe that not only is their performance similar, they appear to have *exactly* the same loss oscillations.

# I  Presentation of Theoretical Results

We model the observed features as a distribution over $x \in \mathbb{R}^{d_1}$, assuming only that its covariance $\Sigma$ exists—for clarity we treat $\Sigma = I$ in the main text. We further model an additional vector $x_o \in \mathbb{R}^{d_2}$ representing the opposing signal, with $d_2 > d_1$. We will suppose that on some small fraction of outliers $p \ll 1$, $x_o \sim \mathrm{Unif}\left(\left\{\pm\sqrt{\frac{\alpha}{pd_2}}\mathbf{1}\right\}\right)$ ($\mathbf{1}$ is the all-ones vector) for some $\alpha$ which governs the feature magnitude, and we let it be $\mathbf{0}$ on the remainder of the dataset. We model the target as the linear function $y = \beta^\top x + \frac{1}{\sqrt{d_2}}\mathbf{1}^\top |x_o|$; this captures the idea that the signal $x_o$ correlates strongly with the target, but in opposing directions of equal strength. Finally, we parameterize the network with vectors $b \in \mathbb{R}^{d_1}, b_o \in \mathbb{R}^{d_2}$ and scalar $c$ in one single vector $\theta$, as $f_\theta(x) = c \cdot (b^\top x + b_o^\top x_o)$. Note the specific distribution of $x_o$ is unimportant—furthermore, in our simulations we observed the exact same pattern with cross-entropy loss. From these experiments and our analysis, it seems that depth and a small signal-to-noise ratio are the only elements needed for this behavior to arise.

A standard initialization would be to sample $[b, b_o]^\top \sim \mathcal{N}(0, \frac{1}{d_1+d_2}I)$, which would then imply highly concentrated distributions for $\|b\|_2^2, \|b_o\|_2^2, b^\top \beta$. As tracking the precise concentration terms would not meaningfully contribute to the analysis, we simplify by directly assuming that at initialization $\|b\|_2^2 = \frac{d_1}{d_1+d_2}$, $\|b_o\|_2^2 = \frac{d_2}{d_1+d_2}$, and $b^\top \beta = \frac{\|\beta\|}{\sqrt{d_1+d_2}}$. Likewise, we let $c = 1$, ensuring that both layers have the same norm. We perform standard linear regression by minimizing the population loss $L(\theta) := \frac{1}{2}\mathbb{E}[(f_\theta(x) - y)^2]$. For our purposes, it will be sufficient to study the trajectory of *gradient flow*—results for gradient descent will then follow for step sizes $\eta$ which are not too large, especially because the timescales of the phases of optimization we study shrink as $\alpha$ grows.

We see that the minimizer of this objective has $b_o = \mathbf{0}$ and $cb = \beta$. However, an analysis of the trajectory of gradient flow will elucidate how depth and large noise lead to sharpening—for larger $\alpha$, the network will be forced initially to minimizes loss on the outliers. We note that in exploring the edge of stability, Cohen et al. [6] found that sometimes the model would have a brief *decrease* in sharpness, particularly for square loss. In fact, we show that this initial decay is expected in the presence of large magnitude, unhelpful signal.

**Theorem I.1** (Initial *decrease* in sharpness). *Let $k := \frac{d_2}{d_1}$, and assume $\|\beta\| > \max\left(\frac{d_1}{\sqrt{d_1+d_2}}, \frac{24}{5}\right)$. At initialization, the sharpness $\|\nabla_\theta^2 L(\theta)\|_2$ lies in $[\alpha, \alpha(1+2\sqrt{\frac{d_2}{d_1+d_2}})]$. Further, if $\sqrt{\alpha} = \Omega(\|\beta\| k \ln k)$, then the sharpness will decrease as $O(e^{-\alpha t})$ from $t = 0$ until some time $t_1 \leq \frac{\ln \|\beta\|/2}{2\|\beta\|}$.*

After this decrease, signal amplification can proceed—but the magnitude of the unhelpful feature means that the sharpness with respect to *how the network uses this feature* will grow, and so a small perturbation to this value will induce a large increase in loss.

**Theorem I.2** (Progressive sharpening). *If $\sqrt{\alpha} = \Omega\left(-\frac{\ln \ln k}{\ln k}\right)$, then at starting at time $t_1$ the sharpness will increase linearly in $\|\beta\|$ until some time $t_2 \geq \frac{1}{2\|\beta\|_2^2}$, reaching at least $\frac{5}{8}\|\beta\|\alpha$.*

The proof of Theorem I.2 shows that the sharpness will occur in $b_o$ in particular, and we observe the same in simulations. Oscillation will not occur during gradient flow—but for SGD with step

28

size $\eta$ larger than $\frac{16}{5\|\beta\|\alpha}$, this result means that $b_o$ will start to increase in magnitude. If this growth continues for long enough, it will rapidly *reintroduce* the outlier feature, which will also cause loss alternation on the opposing outliers. We simulate this model and verify exactly this behavior: in Appendix F we visualize the dynamics alongside an MLP trained on CIFAR-10, which displays the same characteristic behavior. Due to the network being linear, the reintroduction of $x_o$ can only lead to downweighting $x_o$ and returning to the first phase. However, in a non-linear model the sudden reintroduction of a feature which was removed earlier in training seems quite useful for non-linear feature learning (e.g., around iteration 3000 of Fig. 3), where a particular signal may not be useful unless combined with others in a precise way. Though we are unable to reproduce this in our current model, we see the exploration of this behavior and its implications for optimization and generalization as an interesting direction for future study.

## J Proofs of Theoretical Results

Before we begin the analysis, we must identify the quantities of interest during gradient flow and the system of equations that determines how they evolve.

We start by writing out the loss:

$$2L(\theta) = \mathbb{E}[(c(b^\top x + b_o^\top x_o) - (\beta^\top x + d_2^{-1/2} \mathbf{1}^\top |x_o|))^2] \tag{1}$$

$$= \mathbb{E}[((cb - \beta)^\top x)^2] + \mathbb{E}[((cb_o - d_2^{-1/2} \operatorname{sign}(x_o)\mathbf{1})^\top x_o)^2] \tag{2}$$

$$= \|cb - \beta\|^2 + \frac{p}{2} \left( \left( \sqrt{\frac{\alpha}{p}}(cb_o - 1) \right)^2 + \left( \sqrt{\frac{\alpha}{p}}(cb_o + 1) \right)^2 \right) \tag{3}$$

$$= \|cb - \beta\|^2 + \alpha(c^2\|b_o\|^2 + 1). \tag{4}$$

This provides the gradients

$$\nabla_b L = c(cb - \beta), \tag{5}$$

$$\nabla_{b_o} L = \alpha c^2 b_o, \tag{6}$$

$$\nabla_c L = b^\top(cb - \beta) + \alpha\|b_o\|^2 c. \tag{7}$$

We will also make use of the Hessian to identify its top eigenvalue; it is given by

$$\nabla_\theta^2 L(\theta) = \begin{bmatrix} c^2 I_{d_1} & \mathbf{0}_{d_1 \times d_2} & 2cb \\ \mathbf{0}_{d_2 \times d_1} & \alpha c^2 I_{d_2} & 2c\alpha b_o \\ 2cb^\top & 2c\alpha b_o^\top & \|b\|^2 + \alpha\|b_o\|^2 \end{bmatrix}. \tag{8}$$

The maximum eigenvalue $\lambda_{\max}$ at initialization is upper bounded by the maximum row sum of this matrix, and thus $\lambda_{\max} \leq 3\frac{d_1 + \alpha d_2}{d_1 + d_2} < 3\alpha$. Clearly, we also have $\lambda_{\max} \geq \alpha$.

We observe that tracking the precise vectors $b, b_o$ are not necessary to uncover the dynamics when optimizing this loss. First, let us write $b := \epsilon\frac{\beta}{\|\beta\|} + \delta v$, where $v$ is the direction of the rejection of $b$ from $\beta$ (i.e., $\beta^\top v = 0$) and $\delta$ is its norm. Then we have the gradients

$$\nabla_\epsilon L = (\nabla_\epsilon b)^\top (\nabla_b L) \tag{9}$$

$$= \frac{\beta}{\|\beta\|}^\top \left( c^2 \left( \epsilon\frac{\beta}{\|\beta\|} + \delta v \right) - c\beta \right) \tag{10}$$

$$= c^2 \epsilon - c\|\beta\|, \tag{11}$$

$$\nabla_\delta L = (\nabla_\delta b)^\top (\nabla_b L) \tag{12}$$

$$= v^\top \left( c^2 \left( \epsilon\frac{\beta}{\|\beta\|} + \delta v \right) - c\beta \right) \tag{13}$$

$$= c^2 \delta, \tag{14}$$

$$\nabla_c L = \left( \epsilon\frac{\beta}{\|\beta\|} + \delta v \right)^\top \left( c \left( \epsilon\frac{\beta}{\|\beta\|} + \delta v \right) - \beta \right) + \alpha\|b_o\|^2 c \tag{15}$$

$$= c(\epsilon^2 + \delta^2 + \alpha\|b_o\|^2) - \epsilon\|\beta\|. \tag{16}$$

Finally, define the scalar quantity $o := \|b_o\|^2$, noting that $\nabla_o L = 2b_o^\top \nabla_{b_o} L = 2\alpha c^2 o$. Minimizing this loss via gradient flow is therefore characterized by the following ODE on four scalars:

$$\frac{d\epsilon}{dt} = -c^2 \epsilon + c\|\beta\|, \tag{17}$$

$$\frac{d\delta}{dt} = -c^2 \delta, \tag{18}$$

$$\frac{do}{dt} = -2\alpha c^2 o, \tag{19}$$

$$\frac{dc}{dt} = -c(\epsilon^2 + \delta^2 + \alpha o) + \epsilon\|\beta\|. \tag{20}$$

$$\tag{21}$$

30

Furthermore, we have the boundary conditions

$$\epsilon(0) = \sqrt{\frac{1}{d_1 + d_2}}, \tag{22}$$

$$\delta(0) = \sqrt{\frac{d_1 - 1}{d_1 + d_2}}, \tag{23}$$

$$o(0) = \frac{d_2}{d_1 + d_2}, \tag{24}$$

$$c(0) = 1. \tag{25}$$

Given these initializations and dynamics, we make a few observations: (i) all four scalars are initialized at a value greater than 0, and remain greater than 0 at all time steps; (ii) $\delta$ and $o$ will decrease towards 0 monotonically, and $\epsilon$ will increase monotonically until $c\epsilon = \|\beta\|$; (iii) $c$ will be decreasing at initialization. Lastly, we define the quantity $r := (\epsilon(0)^2 + \delta(0)^2 + \alpha o(0)) = \frac{d_1 + \alpha d_2}{d_1 + d_2}$ and $k := \frac{d_2}{d_1}$.

Before we can prove the main results, we present a lemma which serves as a key tool for deriving continuously valid bounds on the scalars we analyze:

**Lemma J.1.** *Consider a vector valued ODE with scalar indices $v_1, v_2, \ldots$, where each index is described over the time interval $[t_{\min}, t_{\max}]$ by the continuous dynamics $\frac{dv_i(t)}{dt} = a_i(v_{-i}(t)) \cdot v_i(t) + b_i(v_{-i}(t))$ with $a_i \leq 0, b_i \geq 0$ for all $i, t$ ($v_{-i}$ denotes the vector $v$ without index $i$). That is, each scalar's gradient is an affine function of that scalar with a negative coefficient. Suppose we define continuous functions $\hat{a}_i, \hat{b}_i : \mathbb{R} \to \mathbb{R}$ such that $\forall i, t, \hat{a}_i(t) \leq a_i(v_{-i}(t))$ and $\hat{b}_i(t) \leq b_i(v_{-i}(t))$. Let $\hat{v}$ be the vector described by these alternate dynamics, with the boundary condition $\hat{v}_i(t_{\min}) = v_i(t_{\min})$ and $v_i(t_{\min}) \geq 0$ for all $i$ (if a solution exists). Then for $t \in [t_{\min}, t_{\max}]$ it holds that*

$$\hat{v}(t) \leq v(t), \tag{26}$$

*elementwise. If $\hat{a}_i, \hat{b}_i$ upper bound $a_i, b_i$, the inequality is reversed.*

*Proof.* Define the vector $w(t) := \hat{v}(t) - v(t)$. This vector has the dynamics

$$\frac{dw_i}{dt} = \frac{d\hat{v}_i}{dt} - \frac{dv_i}{dt} \tag{27}$$

$$= \hat{a}_i(t) \cdot \hat{v}_i(t) + \hat{b}_i(t) - a_i(v_{-i}(t)) \cdot v_i(t) - b_i(v_{-i}(t)) \tag{28}$$

$$\leq \hat{a}_i(t) \cdot \hat{v}_i(t) - a_i(v_{-i}(t)) \cdot v_i(t). \tag{29}$$

The result will follow by showing that $w(t) \leq \mathbf{0}$ for all $t \in [t_{\min}, t_{\max}]$ (this clearly holds at $t_{\min}$). Assume for the sake of contradiction there exists a time $t' \in (t_{\min}, t_{\max}]$ and index $i$ such that $w_i(t') > 0$ (let $i$ be the first such index for which this occurs, breaking ties arbitrarily). By continuity, we can define $t_0 := \max \{t \in [t_{\min}, t'] \; : \; w_i(t) \leq 0\}$. By definition of $t_0$ it holds that $w_i(t_0) = 0$ and $\forall \epsilon > 0, w_i(t_0 + \epsilon) - w_i(t_0) = w_i(t_0 + \epsilon) > 0$, and thus $\frac{dw_i(t_0)}{dt} > 0$. But by the definition of $w$ we also have

$$\hat{v}_i(t_0) = v_i(t_0) + w_i(t_0) \tag{30}$$

$$= v_i(t_0), \tag{31}$$

and therefore

$$\frac{dw_i(t_0)}{dt} \leq \hat{a}_i(t_0) \cdot \hat{v}_i(t_0) - a_i(v_{-i}(t_0)) \cdot v_i(t_0) \tag{32}$$

$$= \big(\hat{a}_i(t_0) - a_i(v_{-i}(t_0))\big) \cdot v_i(t_0) \tag{33}$$

$$\leq 0, \tag{34}$$

with the last inequality following because $\hat{a}_i(t) \leq a_i(v_{-i}(t))$ and $v_i(t) > 0$ for all $i, t \in [t_{\min}, t_{\max}]$. Having proven both $\frac{dw_i(t_0)}{dt} > 0$ and $\frac{dw_i(t_0)}{dt} \leq 0$, we conclude that no such $t'$ can exist. The other direction follows by analogous argument. $\square$

We make use of this lemma repeatedly and its application is clear so we invoke it without direct reference. We are now ready to prove the main results:

## J.1 Proof of Theorem I.1

*Proof.* At initialization, we have $\|\beta\| \geq \frac{d_1}{\sqrt{d_1+d_2}} \implies \|\beta\|\epsilon(0) \geq \frac{d_1}{d_1+d_2} = c(0)(\epsilon(0)^2 + \delta(0)^2)$. Therefore, we can remove these terms from $\frac{dc}{dt}$ at time $t = 0$, noting simple that $\frac{dc}{dt} \geq -\alpha oc$. Further, so long as $c$ is still decreasing (and therefore less than $c(0) = 1$),

$$\frac{d(\|\beta\|\epsilon - c(\epsilon^2 + \delta^2))}{dt} \geq \frac{d(\|\beta\|\epsilon - (\epsilon^2 + \delta^2))}{dt} \tag{35}$$

$$= (\|\beta\| - 2\epsilon)\frac{d\epsilon}{dt} - 2\delta\frac{d\delta}{dt} \tag{36}$$

$$= (\|\beta\| - 2\epsilon)(-c^2\epsilon + \|\beta\|c) - 2\delta(-c^2\delta) \tag{37}$$

$$= -c^2(\epsilon\|\beta\| - 2(\epsilon^2 + \delta^2)) + c(\|\beta\|^2 - 2\epsilon) \tag{38}$$

$$\geq -c(\epsilon\|\beta\| - 2(\epsilon^2 + \delta^2)) + c(\|\beta\|^2 - 2\epsilon) \tag{39}$$

$$= c(\|\beta\|^2 - 2\epsilon - \epsilon\|\beta\| + 2(\epsilon^2 + \delta^2)) \tag{40}$$

$$\geq c(\|\beta\|^2 - \epsilon(2 + \|\beta\|)). \tag{41}$$

Since $c > 0$ at all times, this is non-negative so long as the term in parentheses is non-negative, which holds so long as $\epsilon \leq \frac{\|\beta\|^2}{\|\beta\|+2}$. Further, since $\epsilon c \leq \|\beta\|$ we have

$$\frac{d\epsilon^2}{dt} = 2\epsilon\frac{d\epsilon}{dt} \tag{42}$$

$$= -2c^2\epsilon^2 + 2\epsilon c\|\beta\| \tag{43}$$

$$\leq 2\|\beta\|^2. \tag{44}$$

This implies $\epsilon(t)^2 \leq \epsilon(0)^2 + 2t\|\beta\|^2$. Therefore, for $t \leq \frac{\ln\|\beta\|/2}{2\|\beta\|}$ we have $\epsilon(t)^2 \leq \frac{1}{d_1+d_2} + \|\beta\| \ln\|\beta\|/2 \leq \frac{\|\beta\|^4}{(\|\beta\|+2)^2}$ (this inequality holds for $\|\beta\| \geq 2$). This satisfies the desired upper bound.

Thus the term in Eq. (41) is non-negative for all $t \leq \frac{\ln\|\beta\|/2}{2\|\beta\|}$, and so we have $\frac{dc}{dt} \geq -\alpha oc$ under the above conditions. Since the derivative of $o$ is negative in $c$, a lower bound on $\frac{dc}{dt}$ gives us an upper bound on $\frac{do}{dt}$, which in turn maintains a valid lower bound on $\frac{dc}{dt}$ This allows us to solve for just the ODE given by

$$\frac{dc^2}{dt} = -2\alpha c^2 o, \tag{45}$$

$$\frac{do}{dt} = -2\alpha c^2 o. \tag{46}$$

Define $m := \frac{d_1}{d_1+d_2} = \frac{1}{1+k}$. Recalling the initial values of $c^2, o$, The solution to this system is given by

$$c(t)^2 = \frac{m}{1 - \frac{(1-m)}{\exp(2\alpha mt)}}, \tag{47}$$

$$o(t) = \frac{m}{\frac{\exp(2\alpha mt)}{1-m} - 1} \tag{48}$$

$$= \frac{m}{\exp(2\alpha mt)(1 + k^{-1}) - 1} \tag{49}$$

Since these are bounds on the original problem, we have $c(t)^2 \geq m$ and $o(t)$ shrinks exponentially fast in $t$. In particular, note that under the stated condition $\sqrt{\alpha} \geq \frac{\|\beta\| \ln k}{m(\ln\|\beta\|/2)}$ (recalling $k := \frac{d_2}{d_1} > 1$), we have $\frac{\ln k}{2\sqrt{\alpha}m} \leq \frac{\ln\|\beta\|/2}{2\|\beta\|}$. Therefore we can plug in this value for $t$, implying $o(t) \leq m\left(\frac{d_1}{d_2}\right)^{\sqrt{\alpha}} = mk^{-\sqrt{\alpha}}$ at some time before $t = \frac{\ln\|\beta\|/2}{2\|\beta\|}$.

Now we solve for the time at which $\frac{dc}{dt} \geq 0$. Returning to Eq. (41), we can instead suppose that $\epsilon \leq \frac{\|\beta\|^2 - \gamma}{\|\beta\|+2} \implies \|\beta\|^2 - \epsilon(2 + \|\beta\|) \geq \gamma$ for some $\gamma > 0$. If this quantity was non-negative and

has had a derivative of at least $\gamma$ until time $t = \frac{\ln k}{2\sqrt{\alpha}m}$, then its value at that time must be at least $\frac{\gamma \ln k}{2\sqrt{\alpha}m}$. For $\frac{dc}{dt}$ to be non-negative, we need this to be greater than $c(t)^2 \alpha o(t)$, so it suffices to have

$$\frac{\gamma \ln k}{2\sqrt{\alpha}m} \geq \frac{\alpha m}{\exp(2\alpha m t)(1+k^{-1})-1} \impliedby \gamma \ln k \geq \frac{2\alpha^{3/2}m^2}{\left(\frac{d_2}{d_1}\right)^{\sqrt{\alpha}}(1+k^{-1})-1} \impliedby \gamma \geq \frac{2\alpha^{3/2}m^2 k^{-\sqrt{\alpha}}}{\ln k}.$$ Observe

that the stated lower bound on $\alpha$ directly implies this inequality.

Finally, note that $\|b\|^2 = \epsilon^2 + \delta^2$, and therefore

$$\frac{d\|b\|^2}{dt} = 2\epsilon\frac{d\epsilon}{dt} + 2\delta\frac{d\delta}{dt} \tag{50}$$
$$= -2c^2(\epsilon^2 + \delta^2) + 2c\epsilon\|\beta\|. \tag{51}$$

Since $c(0) = 1$ and $c\epsilon < \|\beta\|$, this means $\|b\|^2$ will also be decreasing at initialization. Thus we have shown that all relevant quantities will decrease towards 0 at initialization, but that by time $t = \frac{\ln k}{2\sqrt{\alpha}m}$, we will have $\frac{dc}{dt} \geq 0$. □

## J.2 Proof of Proof of Theorem I.2

*Proof.* Recall from the previous section that we have shown that at some time $t_1 \leq \frac{\ln k}{2\sqrt{\alpha}m}$, $c(t)^2$ will be greater than $m$ and increasing, and $o(t)$ will be upper bounded by $mk^{-\sqrt{\alpha}}$. Furthermore, $\epsilon(t)^2 \leq \frac{1}{d_1+d_2} + 2t\|\beta\|^2$. To show that the sharpness reaches a particular value, we must demonstrate that $c$ grows large enough before the point $c\epsilon \approx \|\beta\|$ where this growth will rapidly slow. To do this, we study the relative growth of $c$ vs. $\epsilon$.

Recall the derivatives of these two terms:

$$\frac{dc}{dt} = -(\epsilon^2 + \delta^2 + \alpha o^2)c + \|\beta\|\epsilon, \tag{52}$$
$$\frac{d\epsilon}{dt} = -c^2\epsilon + \|\beta\|c. \tag{53}$$

Considering instead their squares,

$$\frac{dc^2}{dt} = 2c\frac{dc}{dt} \tag{54}$$
$$= -2(\epsilon^2 + \delta^2 + \alpha o^2)c^2 + 2\|\beta\|\epsilon c, \tag{55}$$
$$\frac{d\epsilon^2}{dt} = 2\epsilon\frac{d\epsilon}{dt} \tag{56}$$
$$= -2\epsilon^2 c^2 + 2\|\beta\|\epsilon c. \tag{57}$$

Since $\delta, o$ decrease monotonically, we have $\frac{dc^2}{dt} \geq -2(\epsilon^2 + \frac{d_1}{d_1+d_2} + \alpha m\left(\frac{d_1}{d_2}\right)^{\sqrt{\alpha}})c^2 + 2\|\beta\|\epsilon$.

Thus if we can show that $\|\beta\|\epsilon c \geq (\epsilon^2 + 2(\frac{d_1}{d_1+d_2} + \alpha m\left(\frac{d_1}{d_2}\right)^{\sqrt{\alpha}}))c^2$, we can conclude that $\frac{dc^2}{dt} \geq (\epsilon^2 c^2 + \|\beta\|\epsilon c) = \frac{1}{2}\frac{d\epsilon^2}{dt}$—that is, that $c(t)^2$ grows at least half as fast as $\epsilon(t)^2$. And since $\delta, o$ continue to decrease, this inequality will continue to hold thereafter.

Simplifying the above desired inequality, we get

$$\epsilon(\|\beta\| - c\epsilon) \geq c\left(2\frac{d_1}{d_1+d_2} + \alpha mk^{-\sqrt{\alpha}}\right). \tag{58}$$

Noting that $\epsilon(0) \geq \frac{1}{\sqrt{d_1+d_2}}$ at all times, this is implied by

$$\frac{\|\beta\|}{c} - \epsilon \geq 2\frac{d_1}{\sqrt{d_1+d_2}} + \sqrt{d_1+d_2}\alpha mk^{-\sqrt{\alpha}}. \tag{59}$$

33

We also know that $c^2$ has a smaller gradient than $\epsilon^2$ at all timesteps, so $c(t)^2 \leq 1 + 2t\|\beta\|^2$. Further, $\epsilon(t)^2 \leq \epsilon(0) + 2t\|\beta\|^2$. So it is sufficient to show

$$\frac{\|\beta\|}{\sqrt{1 + 2t\|\beta\|^2}} - \sqrt{\frac{1}{\sqrt{d_1 + d_2}} + 2t\|\beta\|^2} \geq 2\frac{d_1}{\sqrt{d_1 + d_2}} + \sqrt{d_1 + d_2}\alpha m k^{-\sqrt{\alpha}} \qquad (60)$$

$$\Longleftarrow \frac{1}{\sqrt{2t}} - \frac{1}{\sqrt[4]{d_1 + d_2}} - \sqrt{2t}\|\beta\| \geq 2\frac{d_1}{\sqrt{d_1 + d_2}} + \sqrt{d_1 + d_2}\alpha m k^{-\sqrt{\alpha}}. \qquad (61)$$

Considering again time $t = \frac{\ln k}{2\sqrt{\alpha}m}$, the LHS becomes $\sqrt{\frac{\sqrt{\alpha}m}{\ln \frac{d_2}{d_1}}} - \frac{1}{\sqrt[4]{d_1 + d_2}} - \sqrt{\frac{\ln \frac{d_2}{d_1}}{\sqrt{\alpha}m}}\|\beta\|$, and plugging in the lower bound on $\sqrt{\alpha}$ shows that Eq. (61) must hold by $t = \frac{\ln k}{2\sqrt{\alpha}m}$, and therefore $\frac{dc^2}{dt} \geq \frac{1}{2}\frac{d\epsilon^2}{dt}$ by some time $t_2 \leq \frac{\ln k}{2\sqrt{\alpha}m}$.

Consider the time $t_2$ at which this first occurs, whereby $c(t_2)^2$ is growing by at least one-half the rate of $\epsilon(t_2)^2$. Here we note that we can derive an upper bound on $c$ and $\epsilon$ at this time using our lemma and the fact that

$$\frac{dc}{dt} \leq \|\beta\|\epsilon, \qquad (62)$$

$$\frac{d\epsilon}{dt} \leq \|\beta\|c. \qquad (63)$$

The solution to this system implies

$$c(t_2) \leq \frac{1}{2}\left(\frac{\exp(\|\beta\|t_2) - \exp(-\|\beta\|t_2)}{\sqrt{d_1 + d_2}} + \exp(\|\beta\|t_2) + 1\right) \qquad (64)$$

$$\leq \frac{1}{2}\left(\exp(\|\beta\|t_2)\left(1 + \frac{1}{\sqrt{d_1 + d_2}}\right) + 1\right) \qquad (65)$$

$$\leq \frac{1}{2}\left(\exp\left(\frac{\|\beta\|\ln k}{2\sqrt{\alpha}m}\right)\left(1 + \frac{1}{\sqrt{d_1 + d_2}}\right) + 1\right), \qquad (66)$$

$$\epsilon(t_2) \leq \frac{1}{2}\left(\exp(\|\beta\|t_2)\left(1 + \frac{1}{\sqrt{d_1 + d_2}}\right) + \frac{1}{\sqrt{d_1 + d_2}} - 1\right) \qquad (67)$$

$$\leq \frac{1}{2}\left(\exp\left(\frac{\|\beta\|\ln k}{2\sqrt{\alpha}m}\right)\left(1 + \frac{1}{\sqrt{d_1 + d_2}}\right) + \frac{1}{\sqrt{d_1 + d_2}} - 1\right) \qquad (68)$$

Then for $\alpha > \left(\frac{\|\beta\|\ln \frac{d_2}{d_1}}{m(\ln \|\beta\| - \ln 2)}\right)^2$, the exponential term is upper bounded by $\frac{\sqrt{\|\beta\|}}{2}$, giving

$$c(t_2) \leq \frac{1}{2}\left(\frac{\sqrt{\|\beta\|}}{2}\left(1 + \frac{1}{\sqrt{d_1 + d_2}}\right) + 1\right) \qquad (69)$$

$$\leq \frac{\sqrt{\|\beta\|}}{2}, \qquad (70)$$

$$\epsilon(t_2) \leq \frac{1}{2}\left(\frac{\sqrt{\|\beta\|}}{2}\left(1 + \frac{1}{\sqrt{d_1 + d_2}}\right) + \frac{1}{\sqrt{d_1 + d_2}} - 1\right) \qquad (71)$$

$$\leq \frac{\sqrt{\|\beta\|}}{2}. \qquad (72)$$

We know that optimization will continue until $\epsilon^2 c^2 = \|\beta\|^2$, and also that $\frac{dc^2}{dt} \geq \frac{1}{2}\frac{d\epsilon^2}{dt}$. Since $c \leq \epsilon$, this implies that $\epsilon^2 \geq \|\beta\|$ before convergence. Suppose that starting from time $t_2$, $\epsilon^2$ grows until

time $t'$ by an additional amount $s$. Then we have

$$s = \epsilon(t')^2 - \epsilon(t_2)^2 \tag{73}$$

$$= \int_{t_2}^{t'} \frac{d\epsilon(t)^2}{dt} \tag{74}$$

$$\leq \int_{t_2}^{t'} 2\frac{dc(t)^2}{dt} \tag{75}$$

$$= 2(c(t')^2 - c(t_2)^2). \tag{76}$$

In other words, $c^2$ must have grown by at least half that amount. Since $\epsilon(t_2)^2 \leq \frac{\|\beta\|}{4}$ and therefore $\epsilon(t')^2 \leq \frac{\|\beta\|}{4} + s$, even if $c(t')^2$ is the minimum possible value of $\frac{s}{2}$ we must have at convergence $\frac{s}{2} = c^2 = \frac{\|\beta\|^2}{\epsilon^2} \geq \frac{\|\beta\|^2}{\frac{\|\beta\|}{4}+s}$. This is a quadratic in $s$ and solving tells us that we must have $s \geq \frac{5}{4}\|\beta\|$. Therefore, $c(t')^2 \geq \frac{5}{8}\|\beta\|$ is guaranteed to occur. Noting our derivation of the loss Hessian, this implies the sharpness must reach at least $\frac{5}{8}\alpha\|\beta\|$ for each dimension of $b_o$. $\qquad\square$

# K Additional Samples Under Various Architectures/Seeds

To demonstrate the robustness of our finding we train a ResNet-18, VGG-11, and a Vision Transformer for 1000 steps with full-batch GD, each with multiple random initializations. For each run, we identify the 24 training examples with the most positive and most negative change in loss from step $i$ to step $i + 1$, for $i \in \{100, 250, 500, 750\}$. We then display these images along with their label (above) and the network's predicted label before and after the gradient step (below). The change in the network's predicted labels display a clear pattern, where certain training samples cause the network to associate an opposing signal with a new class, which the network then overwhelmingly predicts whenever that feature is present.

Consistent with our other experiments, we find that early opposing signals tend to be "simpler", e.g. raw colors, whereas later signals are more nuanced, such as the presence of a particular texture. We also see that the Vision Transformer seems to learn complex features earlier, and that they are less obviously aligned with human perception—this is not surprising since they process inputs in a fundamentally different manner than traditional ConvNets.



(a) Step 100 to 101

(b) Step 250 to 251

(c) Step 500 to 501

(d) Step 750 to 751

Figure 41: **(ResNet-18, seed 1)** Images with the most positive (top 3 rows) and most negative (bottom 3 rows) change to training loss after steps 100, 250, 500, and 750. Each image has the true label (above) and the predicted label before and after the gradient update (below).
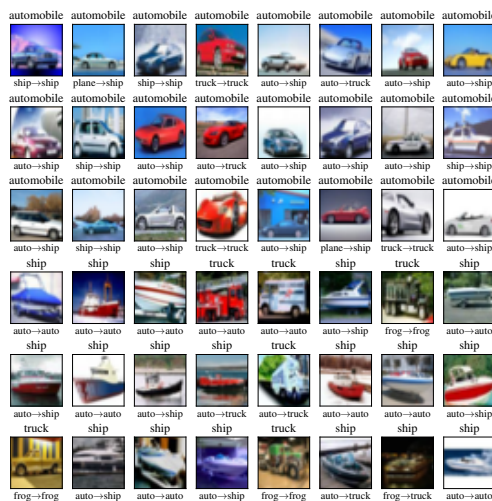
36

(a) Step 100 to 101

(b) Step 250 to 251

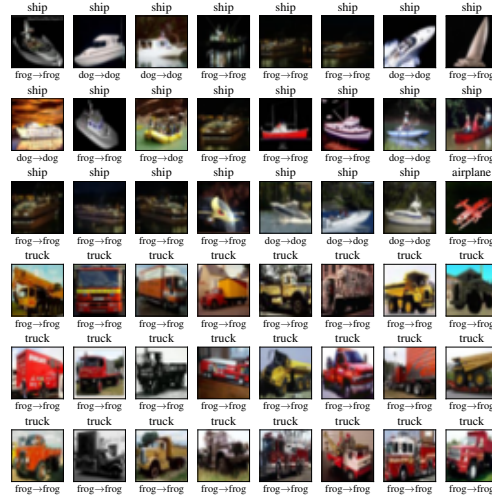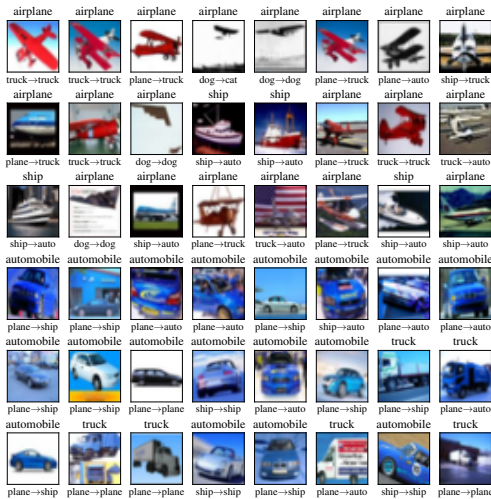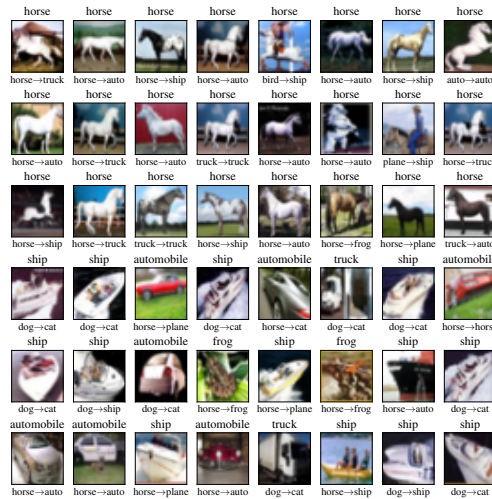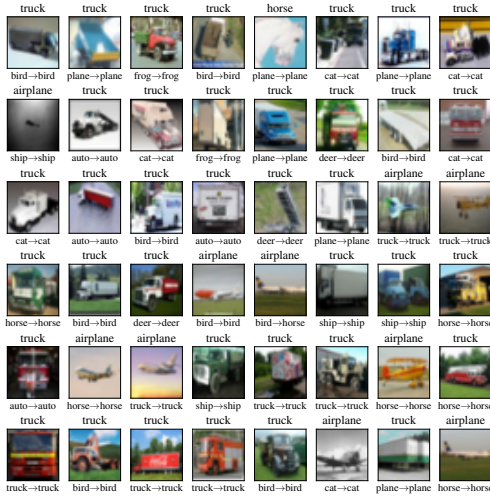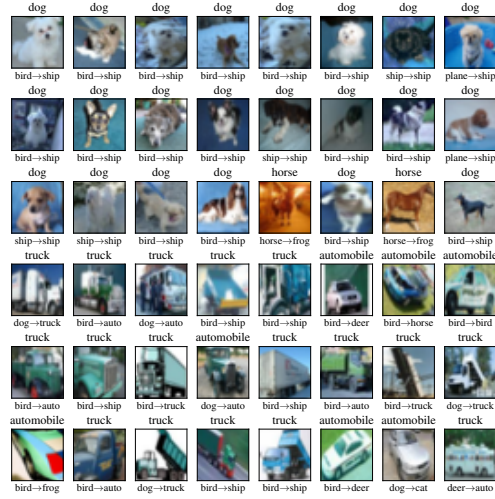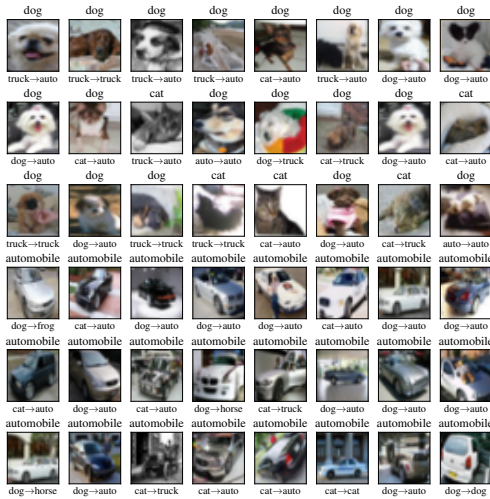(c) Step 500 to 501

(d) Step 750 to 751

Figure 42: **(ResNet-18, seed 2)** Images with the most positive (top 3 rows) and most negative (bottom 3 rows) change to training loss after steps 100, 250, 500, and 750. Each image has the true label (above) and the predicted label before and after the gradient update (below).

(a) Step 100 to 101

(b) Step 250 to 251

(c) Step 500 to 501

(d) Step 750 to 751

Figure 43: **(ResNet-18, seed 3)** Images with the most positive (top 3 rows) and most negative (bottom 3 rows) change to training loss after steps 100, 250, 500, and 750. Each image has the true label (above) and the predicted label before and after the gradient update (below).
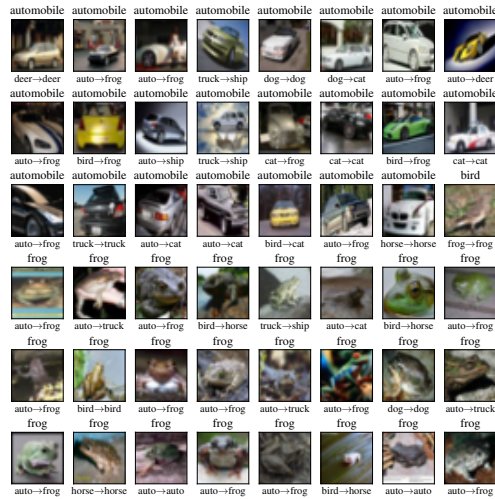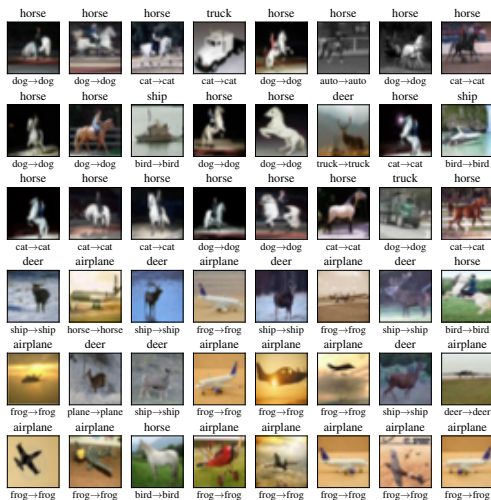
(a) Step 100 to 101

(b) Step 250 to 251

(c) Step 500 to 501

(d) Step 750 to 751

Figure 44: **(VGG-11, seed 1)** Images with the most positive (top 3 rows) and most negative (bottom 3 rows) change to training loss after steps 100, 250, 500, and 750. Each image has the true label (above) and the predicted label before and after the gradient update (below).

(a) Step 100 to 101
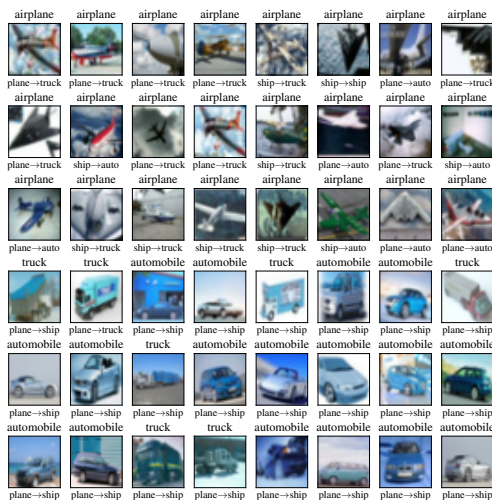
(b) Step 250 to 251
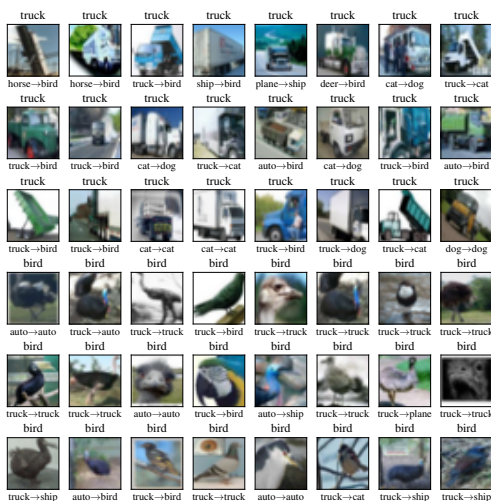
(c) Step 500 to 501

(d) Step 750 to 751

Figure 45: **(VGG-11, seed 2)** Images with the most positive (top 3 rows) and most negative (bottom 3 rows) change to training loss after steps 100, 250, 500, and 750. Each image has the true label (above) and the predicted label before and after the gradient update (below).

(a) Step 100 to 101

(b) Step 250 to 251

(c) Step 500 to 501

(d) Step 750 to 751

Figure 46: **(VGG-11, seed 3)** Images with the most positive (top 3 rows) and most negative (bottom 3 rows) change to training loss after steps 100, 250, 500, and 750. Each image has the true label (above) and the predicted label before and after the gradient update (below).
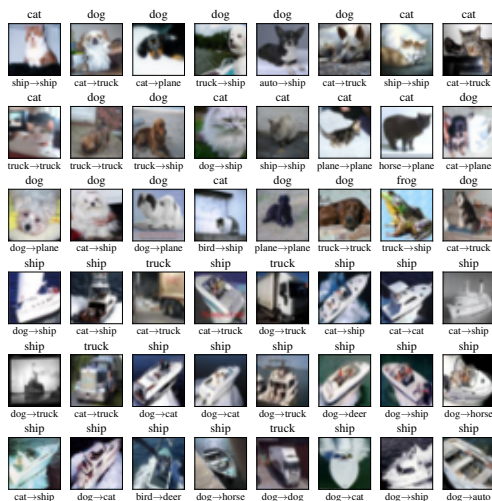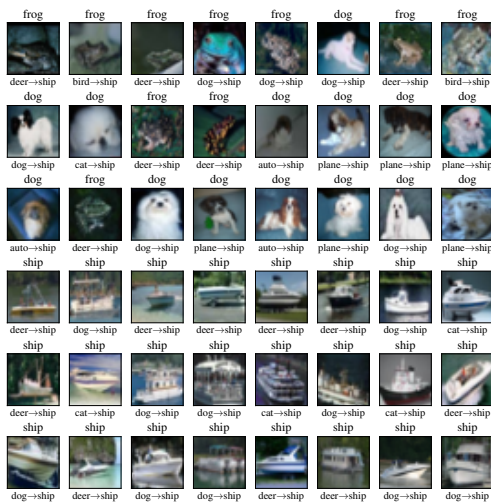
(a) Step 100 to 101

(b) Step 250 to 251

(c) Step 500 to 501

(d) Step 750 to 751

Figure 47: **(ViT, seed 1)** Images with the most positive (top 3 rows) and most negative (bottom 3 rows) change to training loss after steps 100, 250, 500, and 750. Each image has the true label (above) and the predicted label before and after the gradient update (below).

(a) Step 100 to 101

(b) Step 250 to 251

(c) Step 500 to 501
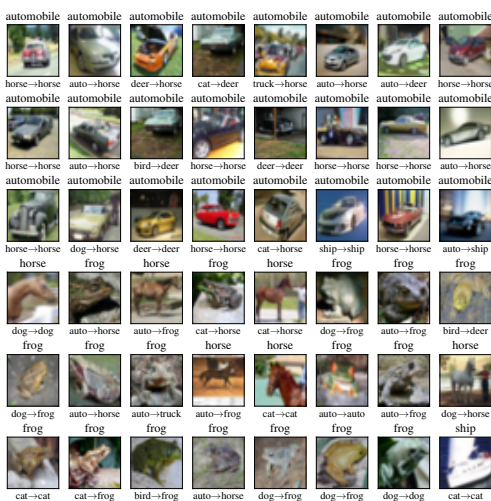
(d) Step 750 to 751

Figure 48: **(ViT, seed 2)** Images with the most positive (top 3 rows) and most negative (bottom 3 rows) change to training loss after steps 100, 250, 500, and 750. Each image has the true label (above) and the predicted label before and after the gradient update (below).
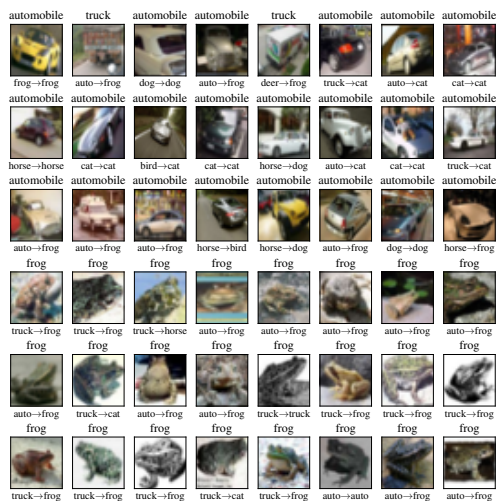
(a) Step 100 to 101

(b) Step 250 to 251

(c) Step 500 to 501

(d) Step 750 to 751

Figure 49: **(ViT, seed 3)** Images with the most positive (top 3 rows) and most negative (bottom 3 rows) change to training loss after steps 100, 250, 500, and 750. Each image has the true label (above) and the predicted label before and after the gradient update (below).